

(19) World Intellectual Property Organization
International Bureau



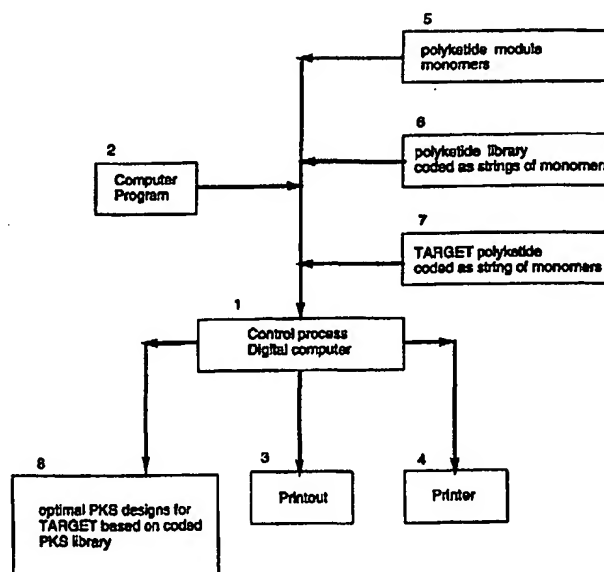
(43) International Publication Date
6 December 2001 (06.12.2001)

PCT

(10) International Publication Number
WO 01/92991 A2

- (51) International Patent Classification⁷: **G06F**
- (21) International Application Number: **PCT/US01/17352**
- (22) International Filing Date: **29 May 2001 (29.05.2001)**
- (25) Filing Language: **English**
- (26) Publication Language: **English**
- (30) Priority Data:
60/207,331 30 May 2000 (30.05.2000) US
60/237,382 4 October 2000 (04.10.2000) US
- (71) Applicant: **KOSAN BIOSCIENCES, INC. [US/US];**
3832 Bay Center Drive, Hayward, CA 94545 (US).
- (72) Inventors: **KHOSLA, Chaitan; 740 La Para Avenue, Palo Alto, CA 94306 (US). REID, Ralph, C.; 600 Galerita Way, San Rafael, CA 94903 (US). SANTI, Daniel, V.; 211 Belgrave Avenue, San Francisco, CA 94417 (US). SIANI, Michael, A.; 341 Day Street, San Francisco, CA 94131 (US).**
- (74) Agents: **FAVORITO, Carolyn, A. et al.; Morrison & Forster LLP, Suite 500, 3811 Valley Centre Drive, San Diego, CA 92130-2332 (US).**
- (81) Designated States (*national*): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, TZ, UA, UG, UZ, VN, YU, ZA, ZW.
- (84) Designated States (*regional*): ARIPO patent (GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE, TR), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).
- Published:
— without international search report and to be republished upon receipt of that report
- For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

(54) Title: **DESIGN OF POLYKETIDE SYNTHASE GENES**



(57) Abstract: Methods for the computational analysis of polyketides and the computer-assisted design of PKS genes are facilitated by representing the structure of a polyketide and/or a PKS gene that encodes the PKS that produces the polyketide by alphanumeric symbols that facilitates computer assisted analysis. A database of polyketides and corresponding PKS genes that can be rapidly searched and information extracted for a variety of applications, including the design and specification of PKS genes via the recombining of modules or portions of modules or sets of modules from already known and variable PKS genes.

DESIGN OF POLYKETIDE SYNTHASE GENES

FIELD OF INVENTION

The present invention provides methods for the analysis of polyketides and the design of polyketide synthase genes. The invention relates to the fields of computational analysis, chemistry, molecular biology, and medicine.

BACKGROUND OF THE INVENTION

The class of compounds known as polyketides is a large family of diverse compounds synthesized primarily from 2-carbon unit building block compounds through a series of condensations and subsequent modifications. Polyketides occur in many types of organisms, including fungi and mycelial bacteria such as the actinomycetes. There are a wide variety of polyketide structures, and the class of polyketides encompasses numerous compounds with diverse activities. Epothilone, erythromycin, FK-506, FK-520, megalomicin, narbomycin, oleandomycin, picromycin, rapamycin, spinocyn, and tylosin are examples of such compounds.

Given the difficulty in producing polyketide compounds by traditional chemical methodology, and the typically low production of polyketides in wild type cells, there has been considerable interest in finding improved or alternate means to produce polyketide compounds. See PCT Publication Nos. WO 95/08548; WO 96/40968; WO 97/02358; and 98/27203; United States Patent Nos. 5,962,290; 5,672,491; and 5,712,146; Fu *et al.*, *Biochemistry* 33: 9321-9326 (1994); McDaniel *et al.*, *Science* 262:1546-1555 (1993); and Rohr, *Angew. Chem. Int. Ed. Engl.* 34(8): 881-888 (1995), each of which is incorporated herein by reference.

Polyketides are synthesized in nature by polyketide synthase (PKS) enzymes. These enzymes, which are complexes of multiple large proteins, are similar to the synthases that catalyze condensation of 2-carbon unit building block compounds in the biosynthesis of fatty acids. The genes that encode PKS enzymes usually consist of three or more open reading frames (ORFs). Two major types of PKS enzymes are known that differ in their composition and mode of synthesis. These two major types of PKS enzymes are commonly referred to as Type I or "modular" and Type II or "iterative" PKS enzymes.

Modular PKSs produce many different polyketides, including a large number of 12-, 14-, and 16-membered macrolide antibiotics including erythromycin, megalomicin, methymycin, narbomycin, oleandomycin, picromycin, and tylosin. Each ORF of a modular PKS can comprise one, two, or more "modules" of ketosynthase activity, each module of which consists of at least two (if a loading module) and more typically three (for the simplest extender module) or more enzymatic activities or "domains." These large multifunctional enzymes (>300,000 kDa) catalyze the biosynthesis of polyketide macrolactones through multistep pathways involving decarboxylative condensations between acyl thioesters followed by cycles of varying β -carbon processing activities (see O'Hagan, D., *The polyketide metabolites*, B. Horwood, New York, 1991, which is incorporated herein by reference).

During the past half decade, the study of modular PKS function and specificity has been greatly facilitated by the plasmid-based *Streptomyces coelicolor* expression system developed with the 6-deoxyerythronolide B (6-dEB) synthase (DEBS) genes (see Kao *et al.*, *Science*, 265: 509-512 (1994), McDaniel *et al.*, *Science* 262: 1546-1557 (1993), and U.S. Patent Nos. 5,672,491 and 5,712,146, each of which is incorporated herein by reference). The advantages to this plasmid-based genetic system for DEBS are that it overcomes the tedious and limited techniques for manipulating the natural DEBS host organism, *Saccharopolyspora erythraea*, allows more facile construction of recombinant PKSs, and reduces the complexity of PKS analysis by providing a "clean" host background. This system also expedited construction of a combinatorial modular polyketide library in *Streptomyces* (see PCT publication No. WO 98/49315, incorporated herein by reference).

The ability to control aspects of polyketide biosynthesis, such as monomer selection and degree of β -carbon processing, by genetic manipulation of PKSs has stimulated great interest in the combinatorial engineering of novel antibiotics (see Hutchinson, *Curr. Opin. Microbiol.* 1: 319-329 (1998); Carreras and Santi, *Curr. Opin. Biotech.* 9: 403-411 (1998); and U.S. Patent Nos. 5,962,290; 5,712,146; and 5,672,491, each of which is incorporated herein by reference). This interest has resulted in the cloning, analysis, and manipulation by recombinant DNA technology of genes that encode PKS enzymes. The resulting technology allows one to manipulate a known PKS gene cluster either to produce the polyketide synthesized by that PKS at higher levels than occur in nature or in hosts that otherwise do not produce the polyketide. The technology also allows one to produce molecules that are structurally related to, but distinct from, the polyketides produced from known PKS gene clusters.

Polyketides are assembled by polyketide synthases through successive condensations of activated coenzyme-A thioester monomers derived from small organic acids such as acetate, propionate, and butyrate. Active sites required for condensation include an acyltransferase (AT), acyl carrier protein (ACP), and beta-ketoacylsynthase (KS). Each condensation cycle results in a β -keto group that undergoes all, some, or none of a series of processing activities. Active sites that perform these reactions include a ketoreductase (KR), dehydratase (DH), and enoylreductase (ER). Thus, the absence of any beta-keto processing domain results in the presence of a ketone, a KR alone gives rise to a hydroxyl, a KR and DH result in an alkene, while a KR, DH, and ER combination leads to complete reduction to an alkane. After assembly of the polyketide chain, the molecule typically undergoes cyclization(s) and post-PKS modification (e.g. glycosylation, oxidation, acylation) to achieve the final active compound.

To illustrate the synthesis of a macrolide by a modular PKS (see Cane *et al.*, *Science* 282: 63 (1998), incorporated herein by reference), one can refer to the PKS that produces the erythromycin polyketide (6-deoxyerythronolide B synthase or DEBS; see U.S. Patent No. 5,824,513, incorporated herein by reference). In the modular DEBS PKS enzyme, the enzymatic steps for each round of condensation and reduction are encoded within a single "module" of the polypeptide (i.e., one distinct module for every condensation cycle). As shown in Figure 1, DEBS consists of a loading module and 6 extender modules and a chain terminating thioesterase (TE) domain within three extremely large polypeptides encoded by three open reading frames (ORFs, designated *eryAI*, *eryAII*, and *eryAIII*).

Each of the three polypeptide subunits of DEBS (DEB1, DEBS2, and DEBS3 in Figure 1) contains 2 extender modules. DEBS1 additionally contains the loading module, and DEBS3 contains the TE domain. Collectively, these proteins catalyze the condensation and appropriate reduction of one propionyl CoA starter unit and six methylmalonyl CoA extender units. Modules 1, 2, 5, and 6 contain KR domains; module 4 contains a complete set, KR/DH/ER, of reductive and dehydratase domains; and module 3 contains no functional reductive domain. Following the condensation and appropriate dehydration and reduction reactions, the enzyme bound intermediate is lactonized by the TE at the end of extender module 6 to form 6-dEB (compound 1 in Figure 1).

More particularly, the loading module of DEBS consists of two domains, an acyltransferase (AT) domain and an acyl carrier protein (ACP) domain. In other PKS enzymes, the loading module is not composed of an AT and an ACP but instead utilizes a partially inactivated

KS, an AT, and an ACP. This partially inactivated KS is in most instances called KS^Q, where the superscript letter is the abbreviation for the amino acid, glutamine, that is present instead of a cysteine in the active site that is believed to be required for condensation activity. Although the KS^Q domain lacks condensation activity, it retains decarboxylase activity. The AT domain of the loading module recognizes a particular acyl-CoA (propionyl for DEBS, which can also accept acetyl) and transfers it as a thiol ester to the ACP of the loading module. Concurrently, the AT on each of the extender modules recognizes a particular extender-CoA (methylmalonyl for DEBS) and transfers it to the ACP of that module to form a thioester. Once the PKS is primed with acyl- and malonyl-ACPs, the acyl group of the loading module migrates to form a thiol ester (trans-esterification) at the KS of the first extender module; at this stage, extender module 1 possesses an acyl-KS and a methylmalonyl ACP. The acyl group derived from the loading module is then covalently attached to the alpha-carbon of the malonyl group to form a carbon-carbon bond, driven by concomitant decarboxylation, and generating a new acyl-ACP that has a backbone two carbons longer than the loading unit (elongation or extension). The growing polyketide chain (various intermediates are shown in Figure 1) is transferred from the ACP to the KS of the next module, and the process continues.

The polyketide chain, growing by two carbons each module, is sequentially passed as a covalently bound thiol ester from module to module, in an assembly line-like process. The carbon chain produced by this process alone would possess a ketone at every other carbon atom, producing a polyketone, from which the name polyketide arises. Commonly, however, additional enzymatic activities modify the beta keto group of the polyketide chain to which the two carbon unit has been added before the chain is transferred to the next module. Modules may contain additional enzymatic activities as well, such as methyl transferase domains, but there are no such additional activities in DEBS.

Once a polyketide chain traverses the final extender module of a modular PKS, it encounters the releasing domain or thioesterase found at the carboxyl end of most PKSs. Here, the polyketide is cleaved from the enzyme and cyclized. The resulting polyketide can be modified further by tailoring or modification enzymes; these enzymes add carbohydrate groups or methyl groups, or make other modifications, i.e., oxidation or reduction, on the polyketide core molecule. For example, the final steps in conversion of 6-dEB to erythromycin A include the actions of a number of modification enzymes, such as: C-6 hydroxylation, attachment of mycarose and desosamine sugars, C-12 hydroxylation (which produces erythromycin C), and

conversion of mycarose to cladinose via O-methylation. These modifications in various combinations result in erythromycins A (compound 2 in Figure 1), B, C, and D.

While the detailed understanding of the mechanisms by which PKS enzymes function and the development of methods for manipulating PKS genes have facilitated the creation of novel polyketides, there remain substantial impediments to the creation of novel polyketides by genetic engineering. One such impediment is the availability of PKS genes. Many polyketides are known but only a relatively small portion of the corresponding PKS genes have been cloned and are available for manipulation. Moreover, in many instances the producing organism for an interesting polyketide is obtainable only with great difficulty and expense, and techniques for its growth in the laboratory and production of the polyketide it produces are unknown or difficult or time-consuming to practice. Also, even if the PKS genes for a desired polyketide have been cloned, those genes may not serve to drive the level of production desired in a particular host cell.

If there were a method to produce a desired polyketide without having to access the genes that encode the PKS that produces the polyketide, then many of these difficulties could be ameliorated or avoided altogether. The present invention meets this need.

SUMMARY OF THE INVENTION

In one embodiment, the present invention provides methods for the computational analysis of polyketides and the computer-assisted design of PKS genes.

In a first aspect, the present invention provides a method for representing the structure of a polyketide and/or a PKS gene that encodes the PKS that produces the polyketide by alphanumeric symbols that facilitates computer assisted analysis.

In a second aspect, the present invention provides a database of polyketides and corresponding PKS genes that can be rapidly searched and information extracted for a variety of applications. More particularly, this database can include, in one mode, all known polyketides; and in another mode, the polyketides, optionally including all intermediates, produced by all known PKS genes or a subset thereof.

In a third aspect, the present invention provides a method for predicting the structure of a PKS and its corresponding genes from the structure of a polyketide.

In a fourth aspect, the present invention provides a method for designing novel PKS genes capable of producing a desired polyketide. This aspect of the invention is directed to the design and specification of PKS genes via the recombining of modules or portions of modules or sets of modules from already known and available PKS genes. In one mode, all possible PKS genes encoding a desired polyketide from a set of genes in a database are generated. In another mode, only a subset of such possible PKS genes is generated based on one or more parameters selected by the user. More particularly, a rating system is provided to sort the PKS genes designed for a particular target polyketide based on any one or more of several criteria, including number of non-native module interfaces, number of non-native protein interfaces, and other parameters as more particularly described below or selected by the user.

In another embodiment, the present invention provides methods and reagents for preparing novel PKS genes that encode PKS enzymes that produce a desired polyketide.

In a first aspect, the present invention provides a library of recombinant DNA compounds, wherein each member of said library encodes a module of a PKS or portions of modules or sets of modules having a desired specificity, and the library as a whole encompasses all of the members of a desired class of specificities.

In a second aspect, the present invention provides a method for assembling a PKS gene cluster that encodes a PKS that produces a desired polyketide from known and available PKS genes other than the naturally occurring PKS genes that produce the polyketide in nature.

These and other embodiments, modes, and aspects of the invention are described in more detail in the following description, the examples, and claims set forth below.

BRIEF DESCRIPTION OF THE FIGURES

Figure 1 shows a schematic representation of the PKS enzyme that synthesizes 6-deoxyerythronolide B (6-dEB, compound 1). The PKS is composed of three proteins, DEBS1, DEBS2, and DEBS3, each of which is represented by an arrow and contains two or more modules. Each module is represented by a solid line, and the domains in each module are shown

inside the arrow. Various intermediates produced during the synthesis are also shown, as are the structures of erythromycins A (compound 2), B, and D resulting from modification of 6-dEB.

Figure 2 shows an illustrative set of 2-carbon unit monomers present in macrocyclic polyketides; these monomers can be used to represent polyketide backbone diversity generated by commonly used starter and extender units (malonyl CoA and methylmalonyl CoA) and the condensation and reduction reactions mediated by PKS enzymes.

Figure 3 shows a representation of 6-dEB by molecular graph, CHUCKLES notation, and SMILES notation. The CHUCKLES notation uses the 2-carbon unit monomers shown in Figure 2. In the CHUCKLES notation, the order of attachment of monomers is designated by the order in which monomers are listed, and the attachment points within the monomers are specified in their definitions. In the SMILES notation, adjacent monomers are attached via single (covalent) bonds depicted by dashes. The cyclization bond is represented by the index 1 adjacent to the Start and Close monomers.

Figure 4 is a flowchart and block flow diagram in five parts designated A-E, inclusive.

Flowchart Figure 4A is a block flow diagram of a computer system to design a novel PKS (and corresponding genes).

Flowchart Figure 4B is a block flow diagram wherein the "Computer Program" block (2) of Flowchart Figure 4A is further defined.

Flowchart Figure 4C is a block flow diagram wherein the "Design novel hybrid PKS genes from library for TARGET" block of Flowchart Figure 4B is further defined.

Flowchart Figure 4D is a block flow diagram wherein the "align TARGET with STARTER; copy to ALIGNMENT" block of Flowchart Figure 4C is further defined.

Flowchart Figure 4E is a block flow diagram wherein the "Rate novel hybrid designs" block (3) of Flowchart Figure 4B is further defined.

Figure 5 shows a flowchart of a matching method for the generation of the CHUCKLES strings used for all polyketides in a library.

DETAILED DESCRIPTION OF THE INVENTION

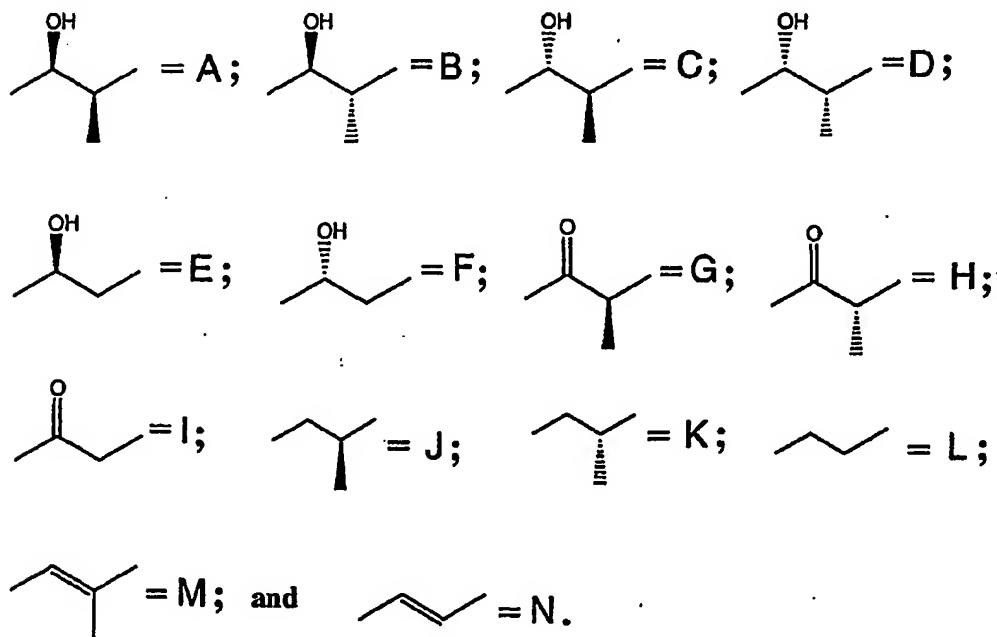
Because polyketides synthesized by modular PKS genes are built by the enzymatically controlled addition of primarily 2-carbon unit monomers and, to a lesser extent, other more complex monomers, each polyketide may be represented as a string of 2-carbon unit and other monomers. These monomers represent the portion of the polyketide backbone structure as a result of the incorporation of various starter and extender units (malonate, methyl malonate, etc.) and the subsequent chemical reactions.

These reactions include:

(1) condensation reactions, of which there are three basic reactions: malonyl-CoA condensation and methylmalonyl-CoA condensation with the branched methyl having either R or S stereochemistry; and

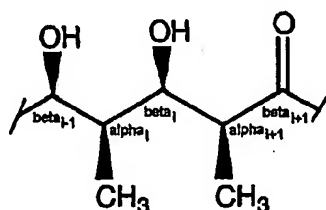
(2) reduction reactions, of which there are five basic reactions: no reduction (ketone preserved), keto-reduction (to yield a hydroxyl having either R or S stereochemistry), dehydration (trans double bond), and enoyl-reduction (to yield a methylene).

An illustrative set of the basic monomers that can be used to represent a polyketide structure (and their corresponding symbols) comprises:

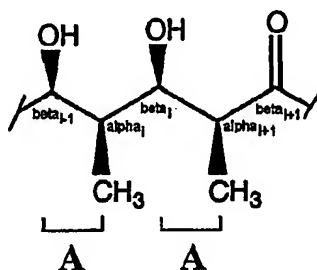


A miscellaneous monomer, Q, can be used to denote a portion of the polyketide structure that cannot be assigned by monomers A-N.

The monomer set shown above and in Figure 2 does not represent the actual monomers incorporated during biosynthesis. Instead, these monomers include a carbon from two different biosynthetic monomers. This is best explained using a polyketide fragment depicted below.



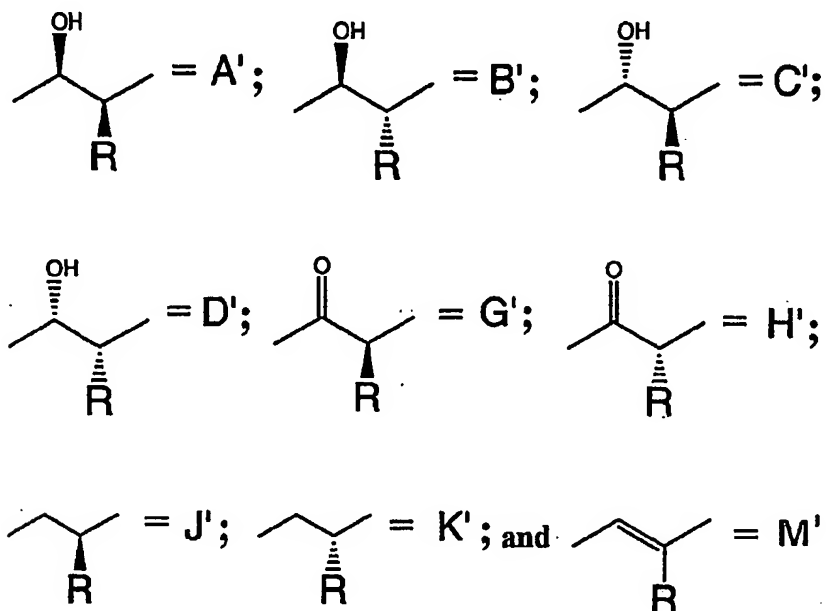
The fragment includes two two-carbon units, i and $i+1$ and part of a third two-carbon unit, $i-1$ that were incorporated into the polyketide during biosynthesis. The i -th extender module attaches the two carbon biosynthetic unit whose backbone carbons are designated as α_i and β_i and the second extender module attaches the two carbon biosynthetic unit whose backbone carbons are designated as α_{i+1} and β_{i+1} . Using the monomer set shown above, this fragment consists of monomer A (derived from the β carbon added in module $i+1$ and the α carbon added in module i) and another monomer A (derived from the β carbon added in module i and the α carbon added in module $i+1$).



The fifth carbon designated β_{i+1} remains unassigned and will depend on the identity of the two-carbon biosynthetic unit that is incorporated in the polyketide by module $i+2$.

The set of monomers shown in Figure 2 can be expanded to include other starter and extender units, of which there are many. Such starter and extender units include, for example but without limitation, hydroxymalonate (e.g., niddamycin), methoxymalonate (e.g. FK-520), ethylmalonate (e.g., FK-520), amino acids or amino acid derivatives that are incorporated into polyketides by the action of a non-ribosomal peptide synthase (e.g., thiazole in epothilone and

pipecolate in rapamycin), or other units incorporated by, for example, an AMP ligase (e.g., the dihydroxycyclohexyl moiety in rapamycin, FK-506, and FK-520) or a soluble CoA ligase. An illustrative set of additional starter and extender units includes:



where R can be anything other than hydrogen or methyl (e.g., allyl, butyl, ethyl, hexyl, hydroxyl, isobutyl, and methoxy).

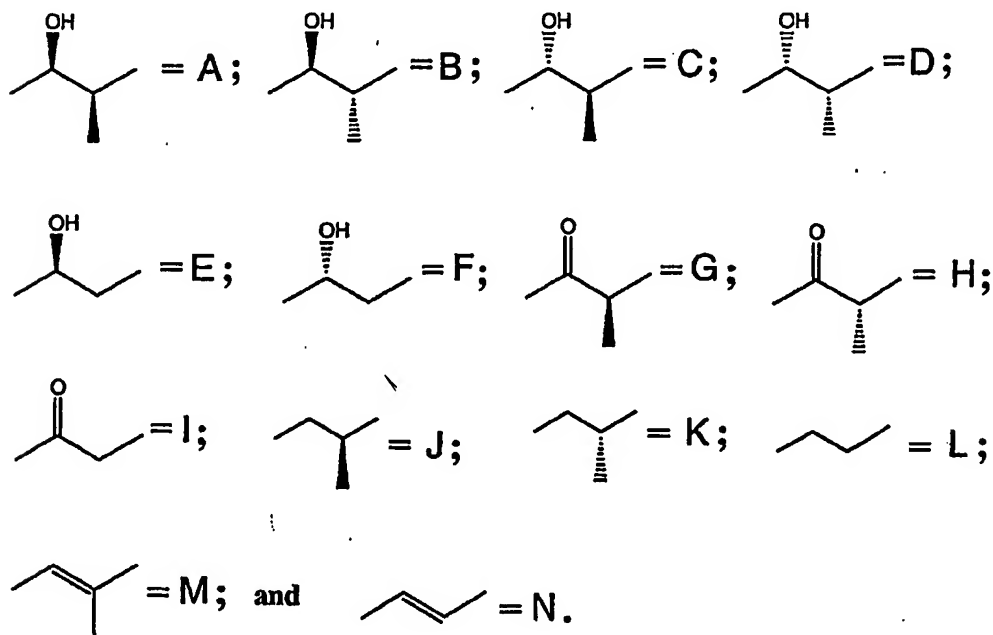
The set of monomers can also include post-PKS modifications, such as hydroxylation, methylation, epoxidation, glycosylation, or addition of intra-macrocyclic fused rings making the system polycyclic. Also, a variety of methods are known for the incorporation of unusual starter and or extender units in polyketide synthases (see, e.g., PCT Publication Nos. WO 97/02358; WO 99/03986; WO 98/01546; and WO 98/01571, each of which is incorporated herein by reference, and the monomer set can include such units.

By viewing polyketides as composed of sets of distinct monomers, one can in accordance with the present invention define a polyketide as a string of alpha-numeric symbols to facilitate computer analysis. In one method, a modified CHUCKLES methodology for representing polyketides is used. The CHUCKLES methodology (see Siani *et al.*, "CHUCKLES: a method for representing and searching peptide and peptoid sequence," *J. Chem. Inf. Sci.* 34: 588-593 (1994) which is incorporated herein by reference) for representing

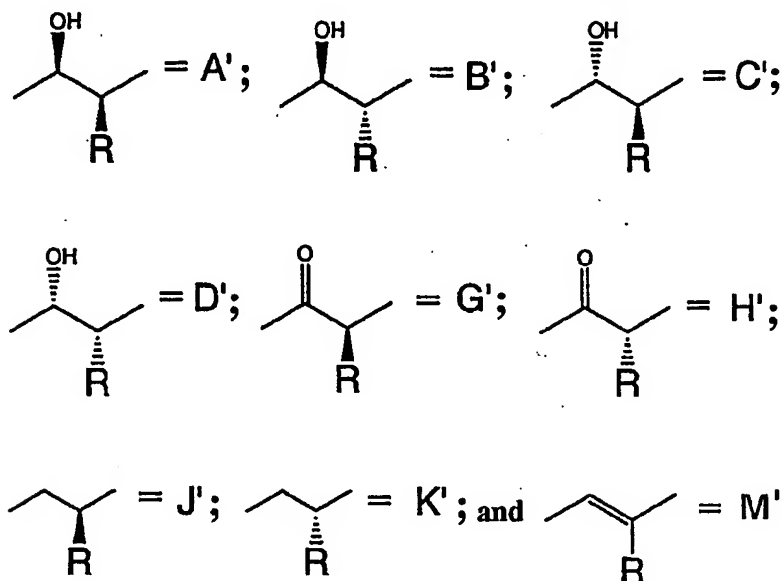
peptides and related oligomers allows monomers to be strung together such that the molecular graph for the basic macrocycle can be generated from the string of monomers.

For example, using the set of monomers comprising A-N described above, the erythromycin macrocycle or 6-dEB can be represented as ADGJDD. This string of alphanumeric symbols is also referred to as the CHUCKLES string. Figure 3 depicts the relationship between the CHUCKLES string, the SMILES string, and the actual molecular structure of 6-dEB. The CHUCKLES string for 6-dEB can be annotated to represent the structure of erythromycin A: A(1-lactone closure,2-hydroxyl)DGJ(2-hydroxyl)D(1-glycosyl)D(1-glycosyl). Thus, ring closure (cyclization) and post-synthetic modifications (glycosylation and hydroxylation), and non-standard units where applicable (there are none in 6-dEB and erythromycin) are entered between parentheses after each monomer. Another example is an annotated CHUCKLES string for epothilone B: ME(1-lactone-closure)M(epoxide)LJDG(2-methylation)E. As above, cyclization, post-synthetic modifications (epoxide formation), and non-standard units (methyl at C-4) are entered between parentheses after each monomer.

In another aspect of the present invention, a database of polyketides is provided. In one aspect of the present invention, the polyketides are represented by a string of defined monomers. In one embodiment, the monomers are selected from a group consisting of:

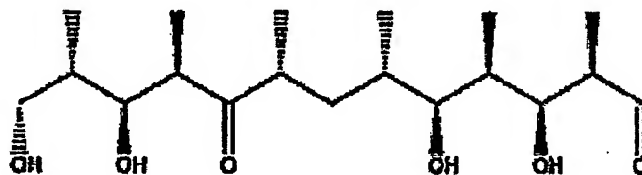


In another embodiment, polyketides are represented by the monomers A-N as well as additional monomers selected from the group consisting of



where R can be anything other than hydrogen or methyl.

The string of monomers can be represented as a linearized structure or as a string of symbols. For example, the erythromycin can be represented as its aglycone, 6-dEB, as



or as a string of symbols, ADGJDD. Optionally, the string of symbols can be annotated as "A(l-lactone closure,2-hydroxyl)DGJ(2-hydroxyl)D(l-glycosyl) D(l-glycosyl)" to more fully capture the erythromycin structure. This set of annotated strings is referred to as a "coded library" or a "coded" database of the present invention.

In an illustrative embodiment, the polyketide database consists of the polyketides described in current literature (Journal of Antibiotics (1981-present), Journal of Natural Products) and various databases (Chemical Abstracts CAPlus, AntiBase). All unique macrocyclic polyketides are converted to the modified CHUCKLES format. Of the ~1000 novel

polyketides obtained, only ~200 different strings of monomers and unique macrocycles are needed to represent the much larger collection of polyketides in the database, because many of the differences between the naturally-occurring polyketides are due to different glycosyl (sugar) groups attached at different positions on the macrocycle.

Thus, a macrocyclic polyketide can be converted to a string of 2-carbon monomers by mapping the monomers onto the polyketide. This can be performed manually or with computer assistance. First, any sugar moieties are conceptually removed by hydrolysis and any lactones (bond between the ketone and oxygen) are hydrolyzed thus generating a linearized structure of the backbone of the polyketide. Generally, this leaves a carboxy carbon at one end of the linear molecule and a hydroxyl at the other. The polyketide is then "sequenced" manually or *in silico* from the end containing the carboxy carbon, the end corresponding to the last monomer added by the PKS before synthesis is complete. This end serves as a convenient handle from which to start the mapping process. Although closing of the lactone often occurs between the two ends of the polyketide, this is not always the case. However, the last ketone added by the PKS is almost always involved in macrolactone formation and so serves as a more convenient handle than the hydroxyl for commencing sequencing.

The manual or *in-silico* sequencing is performed by matching the monomers, one at a time, while traversing the macrocyclic backbone. First the carboxy carbon is skipped, and an attempt is made to match each of the monomers in the monomer set selected (i.e., monomer set A-N in Figure 2) against the next two carbons in the macrocycle. The match takes into account carbon, oxygen, and no substitution at each backbone position, chirality at each backbone position, and bond order between the two backbone carbons.

If the sequencing is performed *in silico*, the method is referred to as back-translation and involves converting a molecular graph into a string of monomers. First, the monomer library is converted to SMARTS format. SMARTS is a superset of the SMILES language that specifies a pattern in a molecular graph (Daylight Software Manual: Theory; Daylight Chemical Information Systems; Irvine, CA 1993, incorporated herein by reference). SMARTS permits one to specify a variable number or a limit on the number of covalent bonds to non-hydrogen atoms from a particular atom. In contrast, SMILES assumes that the unspecified valences are hydrogens. For example, the SMILES string for monomer A is [C@@H](O)[C@H](C). The oxygen may be bonded to any other single atom; if the atom is not specified, it is assumed be a hydrogen. In the SMARTS string for monomer A, [C@@H](O;D2)[C@H]([CH3]), one can

specify the exact number of hydrogens on some atoms (e.g., "CH3"). In addition, the "[O;D2]" indicates the oxygen is bonded to two (from D2) non-hydrogen atoms, in this case the first carbon and some other unspecified atom. This allows matching and distinction of post-modification moieties attached to the oxygen as well as additional cyclizations (six member rings can occur within the macrocycle; e.g., rapamycin). Thus, the SMARTS notation allows pattern matching against the polyketide molecular graph.

When a match occurs, the atoms that match are tagged as part of a superset and labeled with the monomer name. Any atoms that are connected to the monomer that are not part of the macrocycle are tagged for identification as special precursor units (e.g., ethylmalonate instead of methyl malonate or malonate), or post-synthetic modification moieties (e.g., sugars, CCHO, hydroxylation, methylation). If all the atoms and bonds of the monomer cannot be identified, the monomer is given a designation to indicate the lack of identification (e.g., Q for question mark). These Q monomers can be used to identify monomers that are the site of post-PKS modifications that mask the function of the PKS module that generated that portion of the polyketide or that are not in the monomer set and so prevent the correlation of a particular segment of the backbone with one of the monomers in the monomer set.

After a particular 2-carbon unit is identified, the next two carbons are processed the same way. This is repeated until all the backbone carbons are identified and labeled as monomers. When all two-carbon units are identified, one has generated an ordered sequence, or string, of monomers, which is a modified CHUCKLES string of the invention. Moieties corresponding to post-PKS modifications are appended to the monomer in the string as an annotation in parentheses. This method of sequencing may be extended to include any type of monomer. Figure 5 shows a flow chart of this matching method for the generation of the CHUCKLES strings used for all polyketides in a library.

The CHUCKLES string can be in the order corresponding to the direction of biosynthesis on the PKS or its reverse. Each CHUCKLES string has a one-to-many relationship with the PKS gene in the producing organism. Thus, while many different organisms can produce the same polyketide using the same or different PKS genes, each PKS gene generally produces only one PKS that produces only one polyketide (some AT domains can bind different CoAs, leading to the production of multiple polyketides from a single PKS). This allows one to design, from the polyketide structure, a set of PKS genes that would produce that polyketide.

Thus, the present invention provides methods and computational analysis tools for designing PKS genes to produce a desired polyketide. As an illustrative example, the present invention provides a computer program termed MORPH (see the Examples below) that can read the coded library (see the Examples below). An illustrative coded library consists of ~200 unique polyketide CHUCKLES strings. The user specifies the target polyketide, which is converted from molecular structure to a CHUCKLES string.

The program then performs the following, starting with each library compound or string:

- (1) aligns library compound and target compound, emphasizing alignment of adjacent monomers common between the two;
- (2) fills in the gaps using all possible combinations from all library members;
- (3) counts number of non-natural inter-modular boundaries,
- (4) outputs all these alignments.

The alignments are then sorted based on the number of non-natural inter-modular boundaries.

This illustrative program allows one to design and find PKS genes that encode PKS enzymes that are combinations of two or more different PKS enzymes with the fewest inter-modular boundaries, and optionally the fewest inter-protein boundaries. Many other alternative embodiments are provided by the present invention.

For example, one can include the naturally occurring PKS that produces the target polyketide in the coded library to allow components of that PKS to be incorporated into the design of a new PKS. Also, one can include in the coded library non-naturally occurring PKS enzymes, such as those produced and published in the scientific and patent literature to make novel polyketides, in the coded library. See, e.g., PCT publication Nos. WO 98/49315 and WO 96/40968, both of which are incorporated herein by reference.

This CHUCKLES-coded polyketide library can be stored in a computer file as a set of records. In one embodiment, each record contains the chemical name of the polyketide, the unannotated CHUCKLES (containing basic macrocyclic monomers), the annotated CHUCKLES (containing basic macrocyclic monomers with information about post-PKS modifications), the producing organism(s), and other information (e.g., linearized representation

of the polyketide structure, the accession number of organisms or plasmids that have been deposited, gene sequence information, and references).

The MORPH program can read in the polyketide library entries to an array or list of data structures, where each entry data structure contains all or a selected subset of the fields in each library record. The MORPH program then reads in the CHUCKLES-coded TARGET polyketide from the user. This TARGET may optionally be blocked from the library so that it is not used as a STARTER or left in the library, i.e., if it is only distantly related to other known polyketides, or some modules could be useful in designing novel PKS genes, or it is desired to replace only certain PKS modules. This program could also be used for analoging at a particular position via wild-cards defined as part of the TARGET sequence by the user.

Each member of the coded PKS library can be selected as a STARTER unit. Thus, during a run, all library members can be given an equal chance as STARTER units. After a STARTER is chosen, the TARGET is aligned with it. See Flowchart Figure 4D. Any method of alignment can be used such that the maximal number of adjacent STARTER modules is used in the final alignment. After the maximal adjacent modules are used in the ALIGNMENT, smaller adjacent sets or individual modules from the STARTER are used to fill in the gaps. There may be several alignments that are equally good based on the attempt to optimize the number of adjacent modules. For example, if the TARGET contains the "JDG" substring, then 6-dEB, identified with the AID2G3J4D5D6 CHUCKLES string, may align as

| TARGET | JDG |
|--------|--------|
| 6-dEB | J4D2G3 |

or

| TARGET | JDG |
|--------|---------|
| 6-dEB | J4D5G3. |

Both of these alignments have different maximal adjacent modules, with the same length of two (D2G3 in the first and J4D5 in the second). Accordingly, either alignment could be used as STARTERS.

With the optimized alignment from the STARTER, other library entries are systematically used to complete the alignment, or fill in the gaps. This part may be performed on either the optimized ALIGNMENT described above, or the ALIGNMENT without the singl

modules from the STARTER; the removal of the individual modules opens up more space into which larger pieces of the FILLER might be placed. The first library entry is designated as the FILLER. If the FILLER is the same as the STARTER, the next library entry is used as the FILLER. This library entry is flagged as the CURRENT_FILLER_LIBRARY_ENTRY. The same method for finding maximally adjacent modules and then smaller sets or single modules is used to fill the gaps in ALIGNMENT from the FILLER. If not all the gaps are filled in the ALIGNMENT, then the next library entry is used as a new source; that is, it is designated as the FILLER, and the gaps are filled further. This is repeated until the ALIGNMENT is complete or the end of the library is reached.

Assuming all modules in the TARGET are represented in the library, the ALIGNMENT is eventually completely filled. The completed alignment is then written to an output file on the computer disk. When the ALIGNMENT is complete, or there are no more FILLERS in the library, the TARGET and STARTER alignment are re-copied to ALIGNMENT. The CURRENT_FILLER_LIBRARY_ENTRY is incremented, and a new attempt to fill in the gaps is started.

When the CURRENT_FILLER_LIBRARY_ENTRY has reached the end of the library, the ALIGNMENT is wiped, and a new STARTER is chosen. The above process is then repeated for the next STARTER. When all library entries have been used as starters, then all feasible novel polyketide synthases have been generated and written to the computer file. The novel PKSs are then read back into memory and can be further evaluated. An illustrative evaluation process involves:

- (1) counting the non-native inter-module interfaces, and
- (2) counting the number of native inter-protein interfaces (for known and annotated gene sequences).

The novel PKSs are then sorted based on these two numbers, giving higher priority to the non-native inter-module interfaces. In this mode, the goal is to identify those novel PKSs that contain the fewest non-native interfaces.

By providing methods and means for the computer-assisted analysis of polyketides and PKS genes, the present invention greatly facilitates the identification and production of new polyketides with useful activities. Those of skill in the art will appreciate that while the invention is in part illustrated in the Examples below with respect to the design of new PKS

genes for known polyketides, the invention can also be used to design PKS genes for novel polyketides. In this embodiment, one simply provides the structure of the novel polyketide to the MORPH or other program of the invention to generate the desired PKS genes.

Moreover, while the invention is exemplified below by designing new PKS genes composed of the coding sequences for one or more complete modules of two or more different PKS genes, partial modules can also be employed. With the appropriate choice of monomer sets and corresponding coding of the library to be searched, one can generate new PKS gene designs that take advantage of the potential to fuse one PKS gene coding sequence to another at a site corresponding to an intra-modular junction. In another embodiment, one can use "wild-cards" in the encoded polyketide or library to take advantage of known or predicted SAR. Thus, if one knows that a particular position in a polyketide can be varied (i.e., a hydrogen, methyl, or ethyl group at a location determined by an AT domain of a particular module, or a hydroxyl or keto group at a location determined by the presence or absence of a KR domain in a particular module) then one can use a wild-card monomer designation in the polyketide CHUCKLES string to generate PKS genes that produce each of the desired variants.

The methods of the invention have diverse application in addition to the design of new PKS genes. As but one illustrative example, the methods of the invention can be used to design methods to produce a desired compound. Organic molecules containing stereochemical centers are useful for a number of purposes, including use as synthetic or semi-synthetic intermediates. The preparation of such intermediates by organic synthesis can be extremely time consuming and expensive. An alternative source of such intermediates is via specific degradation of a polyketide, and the present invention provides computer-assisted means for designing such production methods.

Thus, certain functional groups of polyketides are susceptible to bond cleavage by specific chemical reactions that do not affect other functional groups. For example, carbon-carbon double bonds can be specifically cleaved by permanganate without affecting other functional groups normally in polyketides, such as ketones, alcohols, and lactones. Likewise, the Baeyer-Villiger reaction converts a ketone to an ester (lactone) without affecting other groups of the aglycone. In accordance with the methods of the invention, one can assemble a library of polyketides in a database that can be addressed with a query describing a particular chemical reaction to generate all of the degradation products produced by that reaction upon each of the polyketides in the library. The degradation fragments thus generated serve as a

library of the invention that can be sorted by properties, such as size, number and type of stereochemical centers, functional groups, or other factors, and searched for useful compounds. Moreover, the functional groups on the ends of the fragments generated (or at other locations) can also be converted to other functional groups by chemical reactions (optionally employing protecting groups on other functional groups), and the database of compounds can be expanded to include the compounds produced by such reactions.

From even a modest library of -200 compounds, one can in this manner generate using the methods of the invention, two to three times as many valuable chemical intermediates. Once such an intermediate is identified, the organism that produces the polyketide from which the fragment is derived is fermented, the polyketide isolated in bulk, the chemical reaction performed, and the desired degradation product(s) isolated and used. In this manner, the present invention makes available a wide variety of useful products otherwise unattainable.

Thus, the present invention has wide application in the fields of chemistry, particularly medicinal chemistry, molecular biology, and medicine. Those of skill in the art will recognize these and other benefits and applications provided by the present invention. Thus, the following examples are given for the purpose of illustrating the present invention and shall not be construed as being a limitation on the scope of the invention or claims.

EXAMPLE 1

The MORPH Program

This example provides the source code for an illustrative MORPH program of the invention. The MORPH program is a command line driven program that runs on a UNIX system. The program can be run from a shell script, such that the user fills in the entire command ahead of time, then post-processes the output file with UNIX utilities including sort, egrep, and uniq.

The command line appears as follows:

```
morph3 -l libraryfile -n targetname -t targetsequence [-x X-wildcards] [-y Y-wildcards] [-z Z-wildcards].
```


The library file is the name of the text file described below in Example 2. The target name is a user-defined identifier to distinguish this target from the library members (e.g., epothiloneD). The target sequence is a string of monomers that represent the CHUCKLES-encoded target polyketide (e.g., MEMLJDGE). Generally, if the target sequence is in the library, it is commented out from the library so that the morph program does not find the target itself. The three different wildcards, X, Y, and Z, are independent sets of monomers that can be included in the target sequence.

The output from the morph program can be redirected to a file. This output file is then post-processed by (1) extracting the HIT lines with valid combinations of modules that yield the target, (2) sorting the HITS based on alphanumeric content using the UNIX sort command, (3) running the UNIX uniq command which removes multiple copies of each HIT, leaving one copy of each, (4) sorting based on the number of pieces in the sequence of modules. Generally, the fewest number of pieces, which correspond to the fewest number of inter-modular interfaces, are desired; these will appear at the top of the output.

Below are some illustrative examples of calls to the MORPH program from a shell script using epothilone as a target. The first example generates combinations that yield epothilone D:

```
%morph3 -l PKS.lib -n epoD -t MEMLJDGE > omorph3_epoD
```

```
%grep HIT omorph3_epoD | sort | uniq | sort +10 -11 > omorph3_epoD.uniq.sort
```

The second example generates combinations that yield a derivative of epothilone D having a hydroxyl at C-13:

```
%morph3 -l PKS.lib -n epoD-13OH -t MEXLJDGE -x ABCD > oepoD-13OH
```

```
%grep HIT oepoD-13OH | sort | uniq | sort +10 -11 > oepoD-13OH.uniq.sort
```

The third example generates an epothilone having wildcards (set 1):

```
%morph3 -l PKS.lib -n epoD-set1 -t MEXYZDgE -x ABCD -y LEFIN -z JACGM > oepoD-set1
```

```
%grep HIT oepoD-set1 | sort | uniq | sort +10 -11 > oepoD-set1.uniq.sort
```

The fourth example generates an epothilone having another set of wildcards (set 2):

```
%morph3 -l PKS.lib -n epoD-set2 -t MEXYZDgE -x JK -y EF -z JACGM > oepoD-set2
```

```
%grep HIT oepoD-set2 | sort | uniq | sort +10 -11 > oepoD-set2.uniq.sort
```

MORPH in its current implementation operates at the monomer level and thus does not handle intra-modular modifications/splitting. Future implementations could convert the

CHUCKLES-encoded strings into the corresponding and equivalent SMILES and then perform more complex chemical analysis of the PKS molecular graphs. Currently, inter-modular double bonds are present in the library, but are ignored by the program. These bonds can introduced post-biosynthetically and the exact source is generally unknown.

The source code for MORPH is found in Appendix A (version 3.0) and B (version 4.0) (deposited in the microfiche appendix).

EXAMPLE 2

Illustrative Polyketide Library

This example provides the contents of an illustrative CHUCKLES encoded polyketide library. The first column provides the name of the polyketide; the second the CHUCKLES string; the third the annotated CHUCKLES string; and the fourth the source organism. Entries under annotated CHUCKLES and source organism are not complete for all of the polyketides.

| POLYKETIDE | CHUCKLES | annotated-CHUCKLES | SOURCE ORGANISM |
|---|------------------------|--|---------------------------|
| 3-acetyl-4"-butyltylosin | FMNGODF | | |
| #aculeximycin | RRNRSMSRSS SSRLSSN | RR(2-ethyl)NRSMSRS(1-glycosyl)RSSS(2-hydroxyl)SR(1-glycosyl)LSSN(2-ethyl) | |
| albocycline-M1-ingramycin-TA2407-cineromycinB-U28010-SR2077 | BLME=JN | BLM(1,2-epoxy)E(1-methoxy)=J(2-hydroxyl)N | |
| albocycline-M2 | BLME=JN | B(2-hydroxyl)LME(1-methoxy)=J(2-hydroxyl)N | |
| albocycline-M3 | BSME=JL | BSME(1-methoxy)=J(2-hydroxyl)L | |
| albocycline-M5 | BSME=JN | BSME(1-methoxy)=J(2-hydroxyl)N | |
| albocycline-M6 | BLME=JL | BL(2-hydroxyl)ME(1-methoxy)=J(2-hydroxyl)L | |
| albocycline-M7 | BLME=JN | BL(2-hydroxyl)ME(1-methoxy)=J(2-hydroxyl)N | |
| albocycline-M8 | BLME=Q | BLME(1-methoxy)=Q | |
| aldgamycin | BMLGJDL | B(1-cyc)MLG(2-hydroxyl)JDL(2-cyc) | |
| amphotericinA | CDNNLNNNNF CEFEALEE | CDNNLNNNNF(1-glycosyl)C(1-O-cyc,2-carboxylicacid)EF(1-cyc)EAL EE | |
| #amphotericinB | CDNNNNNNNF QQQELEE | CDNNNNNNNF(1-glycosyl)C(1-O-cyc,2-carboxylicacid)EF(1-cyc)EAL EE | |
| angiolam | NMFJNSJQLGA m | | |
| aplyronineA | BFJCENFFMEK AFNN | B(1-C(=O)C)F(1-C(=O)C(C)N(C)C)JCENFFME(1-methoxy)KAF(1-C(=O)C(N(C)C)COC)NN | sea hare Aplysia kurodai |
| apoptolidin | EFLMNAMMM | EF(methoxy-1,hydroxy-2)LMNAMMM | |
| aurachinB | MLMLM | MLMLM | |
| aurachinC | MLMLM | MLMLM | |
| A59770 | QQKQQLJFCDN | QQK(2-ethyl)Q=QLJ(2-hydroxyl)F(2-O-glycosyl)CD(2-hydroxyl)N | Amycolatopsis orientalis |
| A82548A-cytovaricin | QQQKQNLJFDD N | | |
| A83543A | FLFQQQ | FLFQQQ | Saccharopolyspora spinosa |
| AB023a | NNNNNRSSSL R | NNNNNRSSSLSR | |
| AH-758 | RSNMURMN | RS(2-methoxy)NMURMN(2-methoxy) | |
| bafilomycinD | BNHCENMKCM N | BNHCE(1-macrocy,2-methoxy)NMKCMN(2-methoxy)Q(keto-macrocy) | S. sp. |
| bafilomycinA1 | BNHCENMKCM N | BNHCE(1-macrocy,2-methoxy)NMKCMN(2-methoxy)Q(keto-macrocy) | |

| POLYKETIDE | CHUCKLES | annotated-CHUCKLES | SOURCE ORGANISM |
|--|------------------------|---|--|
| borrelidin | SNMRUUUS | SNM(2-N)RUUUS | |
| calyculin | QFBBMmNm | QF(1-methoxy)BBMmNm | Discodermia calyx |
| candicidin-candeptin-ascosin-levorin-etc | CDNNNNNNNF CEELIEE | CDNNNNNNNF(1-glycosyl)C(1-O-cyc,2-carboxylicacid)EE(1-cyc)E(2-hydroxyl)LIEE | S. griseus, S. canescens, S. levoris, S. viridoflavus, Stv. grisoviridum |
| candidin | QDNNNNNNNF CEFELIEE | RDNNNNNNNF(1-glycosyl)C(1-O-cyc,2-carboxylicacid)EF(1-cyc)E(2-hydroxyl)LIEE | |
| carbomycin | ENNCODF | EN(1,2-epoxy)NHOF(1-glycosyl,2-methoxy)F(1-C(=O)C) | |
| carbomycinB-magnamycinB | FNNGOFF | FNNGOF(1-glycosyl,2-methoxy)F(1-C(=O)C) | |
| carbomycin-A-magnamycin-deltamycinA4-NSC51001-PS97-3628-WC3628 | ENNHOF | ENNHO(includes CCHO)FF | |
| chalconomycin-myconomycin-aldegamycinDmikonomycin | BNNGKDN | B(2-O-glycosyl)N(1,2-epoxy)NG(2-hydroxyl)KD(1-glycosyl)N | S. bikiniensis, S. albogriseolus |
| chimeramycinB-PTL448 | BMNCODF | B(0-ethyl)MNC(1-glycosyl)OD(1-glycosyl)F | S. ambofaciens ka-448 |
| chivosazola | SRRnNNSRQn RMnNn | SRR(1-O-macrocy)nNNSR(1-methoxy)QnNnR(1-glycosyl)MnNnQ(keto-macrocy) | S. cellulsum |
| cineromycinB | BLME=JN | BLME=J(2-hydroxyl)N | S. cinereochromogena, S. sp. |
| cineromycinBdehydro | BLMI=JN | BLMI=J(2-hydroxyl)N | S. grieviridis |
| cineromycinB2,3dihydro | BLME=JL | BLME=J(2-hydroxyl)L | S. grieviridis |
| ciramycinB1dihydroxy-A6888X | BMNGODF | BM(1,2-epoxy)NGOD(1-glycosyl)F | S. flocculus |
| ciramycinB-ciramycinB1-Acumycin-A688A-B58941-A6888A | AMNGODF | AM(1,2-epoxy)NGOD(1-glycosyl)F | S. cirratius, S. griseoflavus, S. fradiae, S. flocculus |
| cladospolideA | ELLFN | ELLF(2-hydroxyl)N | Cladosporium fulvum, C. cladosporiodes |
| cladospolideB | ELLFn | ELLF(2-hydroxyl)n | Cladosporium fulvum, C. cladosporiodes |
| cladospolideC | ELLEN | ELLE(2-hydroxyl)N | fungus Cladosporium tenuissimum |

| POLYKETIDE | CHUCKLES | annotated-CHUCKLES | SOURCE ORGANISM |
|--|------------------|---|---|
| concanamycinA-folimyacin-A661-1-S45A-TAN1323B-X4357B | CENMKCCMN | CE(2-methoxy)NMKC(2-ethyl)CMN(2-methoxy) | S diastatochromogenese, S. sp, S. neyagawaensis |
| concanamycinB-S45B | CENMKCCMN | CE(2-methoxy)NMKCCMN(2-methoxy) | S diastatochromogenese |
| concanamycinG-anhydroconcanamycinB | NBNHCENMKCMN | NBNHCE(2-methoxy)NMKCCMN(2-methoxy) | |
| conglobatin | AJM | A(1-oxazoyl)JM.AJM | |
| #copiamycin | RNRSSSSSRRLR-SRN | RNRSSSS(1-O-cyc)S(2-hydroxyl)S(1-cyc)RLRSRN | S. hygrosopicus |
| cytovaricin-H230 | QQKQNLJFCDN | QQKQNLJ(2-hydroxyl)F(2-O-glycosyl)CD(2-hydroxyl)N | S. sp., S. collinus |
| cytovaricinB | QQKQNLJFDDN | QQKQNLJ(2-hydroxyl)F(2-O-glycosyl)CD(2-hydroxyl)N | S. torulosus |
| CP64537 | ARGKDD | A(2-glycosyl)R(1-C(=O)C(O)C(C)C)G(2-hydroxyl)KDD(1-glycosyl) | Streptomyces toyocaensis humicola ATCC 39491 |
| damavaricinC | QDQCCNM | QDQCCNM | S. spectabilis |
| deltamycinA1 | ENNGOFF | EN(1,2-epoxy)NGOF(1-glycosyl,2-methoxy)F(1-C(=O)C) | S. deltae, S. halstedii-deltae |
| deltamycinX-desisovalerylcarbo mycinA | ENNGOFF | EN(1,2-epoxy)NGOF(1-glycosyl,2-methoxy)F(1-C(=O)C) | S. deltae, S. halstedii-deltae |
| engleromycin | QNJHN | QNJH(2-hydroxyl)N(1,2-epoxy) | Engleromyces goetzei |
| #epothilone | MEMLJDgE | MEM(1,2-epoxy)LJDG(2-methyl)E | |
| erythromycin | ADGJDD | A(2-hydroxyl)DGJ(2-hydroxyl)D(1-glycosyl)D(1-glycosyl) | |
| espinomycinA2 | ENNCOFF | ENNCOF(1-glycosyl,2-methoxy)F(1-C(=O)C) | S. fungicidicus |
| filipinIII-lagosin14deoxy | ENNNNMEFFFFF | E(2-hydroxyl)NNNNMEFFFFFFFF | S. filipinensis, S. durhamensis |
| filipin-lagosin | ENNNNMEFFFFF | | S. filipinensis, S. durhamensis |
| formamicin | CBNMOCMN | CBNMO(includes a long, branched alkyl chain)CMN | |
| foromacidinB-spiramycinII-spiramycinB | ENNAOFF | | S. ambofaciens |
| FD891 | RRUSNNLNRM | RRUS(1-O-macrocy)NNL(2-hydroxyl)N(1,2-epoxy)RMMQ(keto-macrocy) | S. graminofaciens |
| FK895 | RNRNMRNRLS | R(1-methoxy)N(1,2-epoxy)RNMR(1-O-macrocy)NR(1-C(=O)C,2-hydroxyl)LSQ(keto-macrocy) | S. hygrosopicus |
| FK-506 | MAEPMJTBKQQ | | |
| gedamycin | IEJBNNNNNNNF | IEJBNNNNNNNF(1-glycosyl)A(1-O-cyc,2-carboxylicacid)EF(2-cyc)EEEEIE | S. aureofaciens |

| POLYKETIDE | CHUCKLES | annotated-CHUCKLES | SOURCE ORGANISM |
|---|-----------------------|---|---|
| geldanamycin | QULRMRnM | QUL(2-methoxy)RMS(1-CONH2,2-methoxy)nM | S. hygrosopicus var. gelanus |
| gephyronicacid | RRTSRRM | | |
| gloeosporone | ELLEIL | ELLE(1-O-cyc)I(2-cyc,2-hydroxyl)L | Colletotrichum gloeosporioides f. sp. jussiaea |
| GERI-155 | BNLGJDN | B(2-O-glycosyl)N(1,2-epoxy)LG(2-hydroxyl)JD(1-glycosyl)N | S. GERI-155 |
| halomicin | QNCRCBNM | QNC(1-methoxy)R(1-C(=O)C)CBNM | Mic. halophytica |
| herbimycin | ALDMFnM | A(1-methoxy)L(2-methoxy)D(1-methoxy)MF(1-CONH2,2-methoxy)nM | |
| hygrolidin | CENMJCM | CE(1-O-macrocy,2-methoxy)NMJCMQ(keto-macrocy) | S. hygrosopicus |
| hygrolidin-oxo | BNHCENMJCMN | BNHCE(1-O-macrocy,2-methoxy)NMJCMN(2-methoxy)Q(keto-macrocy) | S. griseus, S. hygrosopicus |
| irumamycin | URRNSLMQQS | UR(1-O-macrocy)RNS(1-glycosyl)LMR(1-O-cyc)=QS(2-cyc)Q(keto-macrocy) | |
| juvenimicinA1-T1124A1-M4365A1 | BMNGFDF | BMNGFDF(includes an ethyl in pos 2) | Mic. chalcea |
| juvenimicinA2-T1124A2-M4365A2 | BMNGJDF | | Mic. chalcea |
| juvenimicinA4-T1124A4-M4365A4 | BMNGODF | | Mic. chalcea-Mic. capillata |
| juvenimicin-T1124 | BNNGJDF | | Mic. chalcea |
| kanchanamycin | NMSSSSSSRL RSRNN | | |
| lankamycin-kujimycin-landavamycin-A20338N2 | ADGJDD | | S. violaceoniger, S. spinichromogenes |
| leinamycin | QNNIMLN | | |
| leucanicidin | CENMKCMN | CE(2-methoxy)NMKCMN(2-methoxy) | S. halstedii |
| leucomycinA12-kitasomycinA12 | FNNCODF | | S. kitasatoensis |
| leucomycinA14-kitasomycinA14 | FnNCOFF | | S. kitasatoensis |
| leucomycinA3-josamycin-platenomycinA3-turimycinA5 | ENNAODF | | S. kitasatoensis, S. hygrosopicus, S. narbonensis, S. platensis |
| leucomycinA5-turimicinH4 | ENNCOFF | | S. kitasatoensis |
| lienomycin | SSSNSSSTSM LNRRNNNNNL | | |

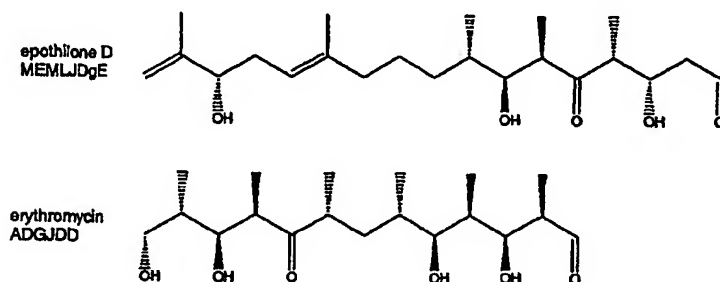
| POLYKETIDE | CHUCKLES | annotated-CHUCKLES | SOURCE ORGANISM |
|--|-----------------------------|-------------------------------|--|
| lucensomycin- etruscomycin- lucimycin-FJ1163- butylpimaricin | SNNNSREEFN N | | Act. sp, S. lucensis, S. glaucus |
| L155175 | RSNMURmM | RS(2-methoxy)NMURmM | |
| L681110 | NMURMN | NMURMN(2-methoxy) | |
| macrocin- lactenocin | BMNHODF | | S. aureus, S. lutea, K. pneum, B. subtl., Shva.. |
| macrocin-YO7625 | QMNGODF | | S. fradiae gs 16 |
| maridomycinL | ENNCOFF | | S. paltensis, S. rimosius, S. capuensis, S. racemochromogene s |
| maridomycin- platenomycinC3- turimycinEP5- B5050A-YL704- C-3 | ENNCJDF | | S. hygroscopicus- S. platensis- malvinus |
| mathemycinA | RSSRSSRSSRR MRMRNLRL | | |
| midecamycinAI- platenomycinB1- SF837 | ENNCOFF | ENNCOF(2-methoxy)F | S. mycarofaciens |
| midecamycinA2- mydecamycinA2- SF837A2 | FNNDJEE | FNNDJE(2-hydroxy)E(1-C(=O)CC) | S. mycarofaciens |
| milbemycin | QQQMKNQQ | | |
| #monazomycin | SSRRSSURNR SMRMRLRL L | | |
| mycinamicin | BNNGJDN | B(1-cyc)NNGJDN(2-cyc) | |
| mycinamicinVI | BNNGJDN | | Micromonospora griseorubida sp. |
| mycinamicinX11 | BNNGLDN | | S. aureus, S. pyogenes, Corvnebacterium |
| #mycolactone | SRMUSMSSL.sS MNMNMN | SRMUSMSSL.sSMNMMN | |
| mycolactoneA | SRMUSMSSL | SRMUSMSSL | |
| mycolactoneB | sSMNMMN | sSMNMMN | |
| myxovirescinA1 | QQEFLNLLIL KJ | | |
| myxovirescinA2 | QQEFLNLLILJ J | | |
| myxovirescinB- mcgovalicinB | QQEFLNLLIL KM | | |
| myxovirescinC-C1 | QQEFLNLLLL KJ | | |
| myxovirescinD | QQEFLNLLLL KM | | |
| myxovirescinE | QQEFLNLLILJ | | |

| POLYKETIDE | CHUCKLES | annotated-CHUCKLES | SOURCE ORGANISM |
|--|--------------------------|-------------------------------|--|
| myxovirescinF1 | QQEFLNNLLLL KJ | | |
| myxovirescinF2 | QQEFLNNLLLL JJ | | |
| myxovirescinG1 | QQEFLNNLLLK J | | |
| myxovirescinG2 | QQEFLNNLLJJ | | |
| myxovirescinH1 | QQEFLNNLLIL KJ | | |
| myxovirescinH2 | QQEFLNNLLIL J | | |
| myxovirescinL | QQEFLNNLLIL HJ | | |
| myxovirescinP1 | QQEFLNNLLIL LKJ | | |
| myxovirescinP2 | QQEFLNNLLIL LJ | | |
| myxovirescinQ | QQEFLNNLLIL KM | | |
| myxovirescinS | QQEFLNNLLIL HJ | | |
| M4365G2 | BMNGODF | | Streptovorticillum kitasatoensis, S. thermaotolerans |
| nancimycin | QNDACBNM | | S. albobovineus |
| neocopiamycin | NRSSSSRLRSR N | | |
| niddamycin- F3463- 3desacetylcarbomy cinB | FNNGOFF | FNNGOF(1-glycosyl,2-methoxy)F | S. aureus, S. lutea, B. subt |
| oligomycinA | QJNNJRGAGAN | | diastatochromogen es, S. chibaensis |
| oligomycinB | QJNNKCHDHD N | | S. diastatochromogen es |
| oligomycinB- 44homo | QJNNJRGRGAN | | S. bottropensis |
| oligomycinD | QJNNJBGAGAN | | S. arabicus, S. parvulus, S. rutgersensis, S. griseus, S. aureofaciens |
| ossamycin | QQQNLLKFDD N | | |
| perimycin | JBNNNNNNNNFC EFEEEEIE | JBNNNNNNNNFCHEEEEEIE | |
| phenalamid | nCNNNNM | | |
| phenalamideA1- fenalamid-I-102-C | JMCMNNNM | | |
| phenalamideA2- 102-T | JMCMNNNM | | |
| phenalamideA3 | JMCMNNNM | | |

| POLYKETIDE | CHUCKLES | annotated-CHUCKLES | SOURCE ORGANISM |
|--|--------------------------------------|------------------------------------|---|
| phenalamideB | JMCMNnNM | | |
| phenalamideC | JMCMNNNM | | |
| phthoramycin | QQNLUSRRN | | |
| pikromycin | ANGJDH | | |
| prasinon-L155175 | DNMJCMM | | S. hydroscopicus ma 5285; S. prasinus |
| protostreptovaricin | QDACDNM | | |
| PD118576A2 | ENMJCMM | | S. sp. wp 3913 |
| PF1163 | EJLLQ | | |
| #quinolidomicin | SSLULSSNNNSI SNNSSSUSRRQ QRRSS | | |
| rapamycin | FGMEGJNNME EKQQ | | |
| rhizopodin | RLQSNSS | | |
| rifamycin | QQQNDACBNM | | |
| rimocidin | BNNNNFCEFEL IA | | |
| rosamicin- repromicin | BMNGODF | | |
| rosaramycin | BMNGODF | | Micromonospora rosaria |
| rustmicin | QQQJMOG | QQQJMO(includes COH)G | |
| rutamycin | QQJNNJBGAGA N | | |
| scytophycin | BFCEENQEMN | | |
| scytophycinB-E | BFCEQNQEMN | | |
| shurimycin | NNRSSSSSRLR SRNN | | |
| sorangicin | LKMFnENLNCF DFNQNNnn | LKMF(2- hydroxyl)nENLNCDFNQNNnn | |
| sorangicinA | QNLNF | | |
| sorangicinB | NLNF | | |
| sorangolideA | LLLLKBLJMEA FM | | myxobacterium sorangium cellulosum |
| soraphen | ELJFDFA | | |
| spiramycin | nNCJDF | | |
| staphcocomycin- angolamycin- shincomycin | CMNGODF | | |
| stipiamide | JMCNNnNM | | |
| tartrolonB2 | nNLFHE | nNLFHE | fragment |
| tedanolide | JGEHMDHF | JGEHMDHF(2-hydroxyl) | |
| thiazinotrienomycin | QLmRSNNNS | | |
| tiacumicin | SMMRMSNM | | |
| tylosinC-macrocin | EMNHODF | | |
| tylosin-A | EMNGODF | | |
| TAN-1323 | NMURRMM | | |
| TMC-34 | NRSSSSSRLRS RN | | |

| POLYKETIDE | CHUCKLES | annotated-CHUCKLES | SOURCE ORGANISM |
|--|-----------|--------------------|-----------------|
| venturicidinA | CKNFLMQQQ | | |
| vicenistatin | LMULRNN | | |
| viranamycinB- virustomycin- TAN1323C | QNMKCCMN | | S. sp. ch41 |
| zincophorin- griseochelin | KMCNLACC | | S. griseus |

In another embodiment, the polyketide library includes the name of the polyketide, the CHUCKLES string and a linearized representation of the structure. The linearized representations of the CHUCKLES structures for erythromycin and epothilone are as follows:

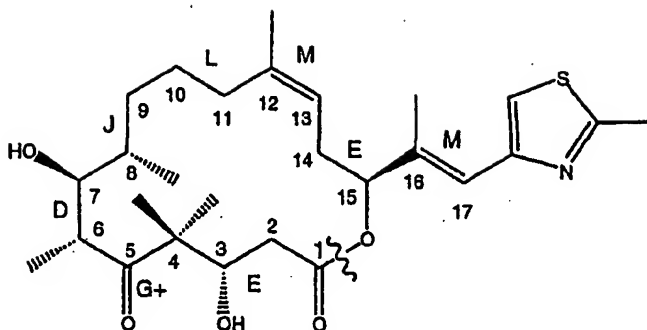


An illustrative example of a polyketide library containing linearized representations of their structures is found in Appendix C (deposited in the microfiche appendix).

EXAMPLE 3

Alternative PKS Genes for Epothilone

This example illustrates the alignment and design of novel PKS genes for the target epothilone. Epothilone is first converted into CHUCKLES string format and then read into the MORPH program as a TARGET. The program then generates all possible alignments of library modules and sorts the alignments to determine preferred combinations of modules for gene construction and production of epothilone via a novel polyketide synthase gene.



The epothilone D structure above was first opened at the macrolactone ring closure between the C-1-ketone and the C-15-oxygen. The monomer set shown in Figure 2 was then matched against each of the successive pairs of macrocyclic backbone carbon atoms, starting with C-2 and C-3, which match monomer E. The next two carbon atoms C-4 and C-5 match monomer G with an additional post-synthetic methylation on C-4. C-6 and C-7 match monomer D. C-8 and C-9 match monomer J. C-10 and C-11 match monomer L. C-12 and C-13 match monomer M. C-14 and C-15, where C-15 has a hydroxyl substitution (modified by thioesterase to close the macrocycle), match monomer E. C-16 and C-17 match monomer M.

The rest of the molecule, a methyl-substituted thiazole moiety, does not match any of the monomers in the monomer set. This moiety corresponds to a malonyl CoA loading module and an NRPS module that together generate the methyl-substituted thiazole moiety. This moiety is thus omitted from the CHUCKLES string generated from this illustrative monomer set but can be added simply by adding a monomer to the set. The CHUCKLES string generated is EGDJLMEM, which is in the reverse order of biosynthesis. This sequence is then reversed to MEMLJDGE to yield a monomer sequence that matches the order of biosynthesis. The sequence is then annotated to account for the post-synthetic modifications as follows MEMLJDG (2-methyl)E.

This target sequence is provided to the MORPH program to generate all possible combinations of modules in the CHUCKLES-encoded library that will yield the target CHUCKLES. The valid combinations are then sorted in increasing order of non-native inter-module interfaces. In one implementation, a MORPH run generated 3,452 valid sequences of five inter-module interfaces. Of these, none contain fewer than five inter-module interfaces. Some illustrative sample module combinations appear below. The combinations are shown listing each monomer followed by a colon and the name of the polyketide(s) from which it is

derived, followed by a parenthetical showing the associated monomers in that polyketide. Vertical lines represent modular junctions between two different polyketides.

Illustrative PKS Gene 1:

M:3acetyl4"-butyryltylosin(FMN) | E:tedanolide(GEH) | M:aldgamycin(BML)
 L:aldgamycin(MLG) | J:aldgamycin(GJD) D:aldgamycin(JDL) | G:tedanolide(JGE)
 E:tedanolide(GEH)

Illustrative PKS Gene 1 thus comprises one or more open reading frames that encode, in the order listed, the module from the acetyl-4"-butyryltylosin PKS that corresponds to monomer M, the module from the tedanolide PKS corresponding to monomer E, the modules from the aldgamycin PKS corresponding to monomers M, L, J, and D, and the modules from the tedanolide PKS corresponding to monomers G and E.

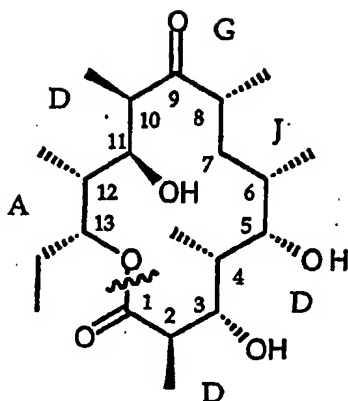
Illustrative PKS Gene 2:

M:albocycline-M1-ingramycin-TA2407-cineromycinB-U28010-SR2077(LME) E:albocycline-M1-ingramycin-TA2407-cineromycinB-U28010-SR2077(MEJ) | M:albocycline-M1-ingramycin-TA2407-cineromycinB-U28010-SR2077(LME) | L:albocycline-M1-ingramycin-TA2407-cineromycinB-U28010-SR2077 (BLM) | J:erythromycin(GJD) D:erythromycin(JDD) | G:tedanolide(JGE) E:tedanolide(GEH)

EXAMPLE 4

Alternative PKS Genes for 6-Deoxvervthrionolide B

This example illustrates the alignment and design of novel PKS genes for the erythromycin basic polyketide structure (6-dEB) using the MORPH program.



For the 6-dEB structure above, the CHUCKLES string is generated by first opening the macrolactone ring closure between the C-1-ketone and the C-13- oxygen. Using the monomer set and matching protocol described in Example 3, one generates the CHUCKLES string DDJGDA, in the reverse order of biosynthesis. This sequence is then reversed to ADGJDD to yield the monomer sequence that matches the order of biosynthesis. The sequence is then annotated to account for the post-synthetic modifications (erythromycin A) as follows A(Z-hydroxyl) DGJ(2-hydroxyl)D(1-glycosyl)D(1-glycosyl).

This target sequence is supplied to the MORPH program to generate all possible combinations of modules in the CHUCKLES-encoded library. The valid combinations are then sorted in increasing order of non-native inter-module interfaces. In one implementation, a MORPH run generated 19,631 valid sequences of less than or equal to five inter-module interfaces. Of these, 13,306 contain 4 inter-module interfaces, and 256 contain only 3 inter-module interfaces. Some of these contain only two inter-module faces, and one only contains one. Some illustrative sample module combinations follow.

Illustrative PKS Gene 1:

A:amphotericinA(EAL) | D:aldgamycin(JDL) | G:mycinamicin(NGJ) J:mycinamicin(GJD)
D:mycinamicin(JDN) | D:amphotericinA(CDN)

Illustrative PKS gene 1 thus comprises one or more open reading frames that encode, in the order listed, the amphotericin PKS module corresponding to monomer A, the aldgamycin PKS monomer corresponding to monomer D, the mycinamicin PKS modules corresponding to monomers G, J, and D, and the amphotericin PKS module corresponding to monomer D.

Illustrative PKS Gene 2:

A:amphotericinA(EAL) | D:aldgamycin(JDL) | G:pikromycin(NGJ) J:pikromycin(GJD)
 D:pikromycin(JDH) | D:aldgamycin(JDL)

Illustrative PKS Gene3:

A:lankamycin-kujimycin-landavamycin-A20338N2(-AD) D:lankamycin-kujimycin-
 landavamycin-A20338N2(ADG) G:lankamycin-kujimycin-landavamycin-A20338N2(DGJ)
 J:lankamycin-kujimycin-landavamycin-A20338N2(GJD) | D:ossamycin(FDD)
 D:ossamycin(DDN)

Illustrative PKS Gene 4:

A:amphotericinA(EAL) | D:lankamycin-kujimycin-landavamycin-A20338N2(ADG)
 G:lankamycin-kujimycin-landavamycin-A20338N2(DGJ) J:lankamycin-kujimycin-
 landavamycin-A20338N2(GJD) | D:A82548A-cytovaricin(FDD) D:A82548A-
 cytovaricin(DDN)

Illustrative PKS Gene 5:

A:lankamycin-kujimycin-landavamycin-A20338N2(-AD) D:lankamycin-kujimycin-
 landavamycin-A20338N2(ADG) G:lankamycin-kujimycin-landavamycin-A20338N2(DGJ)
 J:lankamycin-kujimycin-landavamycin-A20338N2(GJD) D:lankamycin-kujimycin-
 landavamycin-A20338N2(JDD) D:lankamycin-kujimycin-landavamycin-A20338N2(DD-)

EXAMPLE 5

Source Code:

```
#include <stdio.h>
```

```
/* ~siani/programs/morph/morph3.c
```

PURPOSE: To traverse recursively all the entries in PKS.lib,
 generating all feasible combinations of PKS modules to make the TARGET (e.g., epothilone).

INPUT: libraryfile: tab-delimited CHUCKLES-coded polyketides file with the following columns:

1. polyketide name
 2. plain CHUCKLES
 3. annotated CHUCKLES (contains information about post-synthetic modifications)
 4. source organism;
- targetname: user-defined name (e.g., epoD);
- targetsequence: CHUCKLES-coded polyketide of desired TARGET (e.g., MEMLJDGE);

X, Y, Z sets of wildcards: sets of monomers for particular positions appearing in target sequence (the wildcards can be used for analoging the TARGET polyketide);

hard-coded parameters which may be reset (requires recompiling):

NBOUNDARY_CUTOFF determines the maximum number of non-native inter-modular interfaces which are contained in the output (set to 5, but may be increased or decreased); and

RECURSION_COUNTER_CUTOFF specifies the number of levels of recursion (defaults to 0, 1, 2) acceptable for the run -- a large PKS library can cause recursion that will greatly increase run time; because of the multi-directionality of the alignments (using every library entry as a STARTER), there is typically no need to go beyond 2 levels of recursion.

OUTPUT:

All combinations of modules that meet parameters set by user. Example output from MEMLJDGE (epothilone D) using a subset of a PKS library is provided below. Vertical bars indicate non-native inter-modular interfaces. Last column contains the number of "pieces" that are needed to put together the PKS.

Names of PKSs have been abbreviated to fit them in these comments.

```
HIT M:3atyl(FMN)| E:tedan(GEH)| M:aldga(BML) L:aldga(MLG)|
J:aldga(GJD) D:aldga(JDL)| G:tedan(JGE) E:tedan(GBH)| 5
HIT M:albM1(LME) E:albM1(MEJ)| M:albM1(LME)| L:aldga(MLG)|
J:aldga(GJD) D:aldga(JDL)| G:tedan(JGE) E:tedan(GBH)| 5
HIT M:albM1(LME) E:albM1(MEJ)| M:aldga(BML) L:aldga(MLG)|
J:aldga(GJD) D:aldga(JDL)| G:3atyl(NGO)| E:albM1(MEJ)| 5
HIT M:albM1(LME) E:albM1(MEJ)| M:aldga(BML) L:aldga(MLG)|
```

```
J:aldga(GJD) D:aldga(JDL) G:aldga(LGJ) E:albM1(MEJ) 5
HIT M:albM1(LME) E:albM1(MEJ) M:aldga(BML) L:aldga(MLG)
J:aldga(GJD) D:aldga(JDL) G:aldga(LGJ) E:albM1(MEJ) 5
```

USAGE:

```
morph3 -l libraryfile -n targetname -t targetsequence [-x X-wildcards] [-y Y-
wildcards] [-z Z-wildcards]
```

examples:

```
# generate combinations that yield epothilone D
```

```
%morph3 -l PKS.lib -n epoD -t MEMLJDGE > omorph3_epoD
```

```
%egrep HIT omorph3_epoD | sort | uniq | sort +10 -11 >
```

```
omorph3_epoD.uniq.sort
```

```
%egrep ALIGN_TARGET omorph3_epoD >
```

```
omorph3_epoD_STARTER_ALIGN
```

```
# generate combinations that yield epothilone D with a
C13-hydroxyl
```

```
%morph3 -l PKS.lib -n epoD-13OH -t MEXLJDGE -x ABCD >
```

```
oepoD-13QH
```

```
%egrep HIT oepoD-13OH | sort | uniq | sort +10 -11 >
```

```
oepoD-13OH.uniq.sort
```

```
%egrep ALIGN_TARGET oepoD-13OH > oepoD-13OH_STARTER_ALIGN
```

```
# generate combination that yield epothilone with the
following wildcards (set 1)
```

```
%morph3 -l PKS.lib -n epoD-set1 -t MEXYZDgE -x ABCD -y LEFIN -z
```

```
JACGM > oepoD-set1
```

```
%grep HIT oepoD-set1 | sort | uniq | sort +10 -11 > oepoD-set1.uniq.sort
```

```
# generate combination that yield epothilone with the following wildcards (set 2)
```

```
%morph3 -l PKS.lib -n epoD-set2 -t MEXYZDgE -x JK -y EF -z JACGM >
```

```
oepoD-set2
```

```
%grep HIT oepoD-set2 | sort | uniq | sort +10 -11 > oepoD-set2.uniq.sort
```

LIMITATIONS:

This version does not handle intra-modular modifications/splitting because morph is operating at the monomer level. Modifications could convert the CHUCKLES-encoded strings into the corresponding and equivalent SMILES and then perform more complex chemical analysis of the PKS molecular graphs.

Currently, inter-modular double bonds are present in the library, but are ignored by the morph program.

```
*/  
#include <stdio.h>  
/* ~siani/programs/morph/morph3.c  
*/  
#define TRUE 1  
#define FALSE 0  
#define DEBUG_MATCH FALSE  
#define DEBUG_STARTER FALSE  
#define DEBUG_ALIGN FALSE  
#define DEBUG_RECURSE FALSE  
#define DEBUG_WILDCARD FALSE  
#define MAXLEN 80  
#define MAX_TYL_LEN 6  
#define MAX_EPO_LEN 6  
#define MAXNAMELEN 160  
#define MAX_LIB_ENTRIES 500  
#define MAXWILD 3  
#define MAXBUF 200  
#define NBOUNDARY_CUTOFF 5  
#define RECURSION_COUNTER_CUTOFF 2  
#define STARTER_MINIMUM_ADJACENT_ALIGN 2  
#define MINIMUM_ADJACENT_ALIGN 2  
typedef struct _lib {  
    char name[MAXNAMELEN];  
    char monomersequence[MAXNAMELEN];  
    char annotatedsequence[MAXNAMELEN];  
    char alignedsequence[MAXNAMELEN];  
    char alignedPKSname[MAXLEN][MAXNAMELEN];  
    int boundarytoright[MAXNAMELEN];  
    int marked[MAXLEN];  
    char context[MAXLEN][4];  
    int recursion_tagged;  
    int nboundary;
```

```
} LIB;
main(int argc, char **argv)
{
    int    ii=0, jj=0, kk=0, ll=0;
    int    nlib=0;
    int    ecoun=0;
    int    nfilled=0, nfilledmax=0;
    int    epothilonelen=0;
    int    nlargestpiece=0, tnlargestpiece=0;
    int    mmpass=0;
    int    lcount=0;
    int    new_unmarked_entries_filled=0;
    int    recursion_counter = 0;
    int    nwildcard=0;
    int    best_new_unmarked_entries_filled = 0;
    int    smallest_acceptable_piece = 0;
    int    current_nmarked=0, previous_nmarked=0;
    char    *sptr, *eptr, *lptr, *bufptr;
    char    *clibptr;
    char    *libraryfile;
    char    *targetsequence, *targetname;
    char    buf[MAXBUF];
    char    wildcards[MAXWILD][MAXLEN];
    FILE    *libp;
    LIB    epotemp;
    LIB    library[MAX_LIB_ENTRIES];
    LIB    epothilone;
    char    *programe;
    char    **filelist, **fileptr;
    libraryfile = "";
    targetsequence = "";
    targetname = "";
    for(ii=0; ii<MAXWILD; ii++) {
        for(jj=0; jj<MAXLEN; jj++) {
            wildcards[ii][jj] = '\0';
```

```

    }
}
/* process arguments */
filelist = fileptr = (char **)(malloc(argc * sizeof(*argv)));
progname = *argv++;
if(argc<2) {
    fprintf(stderr,"usage:%s -l libraryfile -n targetname -t targetsequence [-x X-
wildcards] [-y Y-wildcards] [-z Z-wildcards]\n",progname);
    exit();
}
while(argc-- > 1) {
    if(argv[0][0] == '-' && argv[0][1] != '\0') {
        /* handle option */
        *++(*argv); /* advance past the minus */
        switch(**argv) {
            case 'l': /* get library input filename (PKS.lib) */
                argv++; argc--;
                libraryfile = argv[0];
                fprintf(stderr,"-l: libraryfile=%s\n",libraryfile);
                break;
            case 'n': /* get target name string */
                argv++; argc--;
                targetname = argv[0];
                fprintf(stderr,"-t: targetname=%s\n",targetname);
                break;
            case 't': /* get target sequence string */
                argv++; argc--;
                targetsequence = argv[0];
                fprintf(stderr,"-t: targetsequence=%s\n",targetsequence);
                break;
            case 'x': /* get a wildcard string */
                argv++; argc--;
                strcpy(wildcards[0],argv[0]);
                fprintf(stderr,"-x:
wildcards[%d]=%s\n",0,wildcards[0]);

```

```
        nwildcard++;
        break;
    case 'y': /* get a wildcard string */
        argv++; argc--;
        strcpy(wildcards[1],argv[0]);
        fprintf(stderr,"-x:
wildcards[%d]=%s\n",1,wildcards[1]);
        nwildcard++;
        break;
    case 'z': /* get a wildcard string */
        argv++; argc--;
        strcpy(wildcards[2],argv[0]);
        fprintf(stderr,"-x:
wildcards[%d]=%s\n",2,wildcards[2]);
        nwildcard++;
        break;
    case 'X': /* get a wildcard string */
        argv++; argc--;
        strcpy(wildcards[0],argv[0]);
        fprintf(stderr,"-x:
wildcards[%d]=%s\n",0,wildcards[0]);
        nwildcard++;
        break;
    case 'Y': /* get a wildcard string */
        argv++; argc--;
        strcpy(wildcards[1],argv[0]);
        fprintf(stderr,"-x:
wildcards[%d]=%s\n",1,wildcards[1]);
        nwildcard++;
        break;
    case 'Z': /* get a wildcard string */
        argv++; argc--;
        strcpy(wildcards[2],argv[0]);
        fprintf(stderr,"-x:
wildcards[%d]=%s\n",2,wildcards[2]);
```

```
        nwildcard++;
        break;
    default:
        fprintf(stderr,"%s unknown option;
ignored\n",*argv);
        /*switch*/
    } else { /* a regular filename */
        *fileptr++ = *argv;
        *fileptr = NULL;
    }
    argv++;
} /*while*/
if(nwildcard > 0) {
    for(ii=0; ii< nwildcard; ii++) {
        fprintf(stderr,"wildcards[%d]=%s\n",ii,wildcards[ii]);
        fprintf(stdout,"wildcards[%d]=%s\n",ii,wildcards[ii]);
    }
}
epothilone.nboundary = 0;
for(ii=0; ii<MAXNAMELEN; ii++) {
    epothilone.name[ii] = '\0';
    epothilone.monomersequence[ii] = '\0';
    epothilone.alignedsequence[ii] = '\0';
    epothilone.boundarytoright[ii] = TRUE;
    for(jj=0; jj<MAXLEN; jj++) {
        epothilone.alignedPKSname[jj][ii] = '\0';
        epothilone.marked[jj] = FALSE;
        epothilone.context[jj][0] = '\0';
        epothilone.context[jj][1] = '\0';
        epothilone.context[jj][2] = '\0';
    }
}
strcpy(epothilone.name,targetname);
strcpy(epothilone.monomersequence,targetsequence);
fprintf(stdout, "TARGET:   %s\n", epothilone.monomersequence);
```

```

        ecount = 0;
        eptr = epothilone.monomersequence;
        while(*eptr != '\0'){
            if(ecount == 0) {
                epothilone.context[ecount][0] = '-';
                epothilone.context[ecount][1] = *eptr;
                epothilone.context[ecount][2] = *(eptr + 1);
            } else {
                if(ecount ==
(strlen(epothilone.monomersequence) - 1)){
                    epothilone.context[ecount][0] =
*(eptr - 1);
                    epothilone.context[ecount][1] = *eptr;
                    epothilone.context[ecount][2] = '-';
                } else {
                    epothilone.context[ecount][0] =
*(eptr - 1);
                    epothilone.context[ecount][1] = *eptr;
                    epothilone.context[ecount][2] =
*(eptr + 1);
                }
            }
            epothilone.context[ecount][3] = '\0';
            eptr++;
            ecount++;
        }
        for(ii=0; ii<ecount; ii++) {
            fprintf(stdout, "(%s)\n", epothilone.context[ii]);
        }
        /* ----- library ----- */
        nlib = get_library(libraryfile, library);
        fprintf(stdout, "nlib=%d\n", nlib);
        kk=0; while(kk < nlib) {
            /* ----- zero out the epothilone entry with respect to a new alignment ----- */
            for(ii=0; ii<MAXNAMELEN; ii++) {

```

```

        epothilone.alignedsequence[ii] = '\0';
        epothilone.boundarytoright[ii] = TRUE;
        for(jj=0; jj<MAXLEN; jj++) {
            epothilone.alignedPKSname[jj][ii] = '\0';
            epothilone.marked[jj] = FALSE;
        }
    }
    /* - reset the context back to that in epothilone- */
    ecount = 0;
    eptr = epothilone.monomersequence;
    while(*eptr != '\0'){
        if(ecount == 0) {
            epothilone.context[ecount][0] = '-';
            epothilone.context[ecount][1] = *eptr;
            epothilone.context[ecount][2] = *(eptr + 1);
        } else {
            if(ecount ==
(strlen(epothilone.monomersequence) - 1)){
                epothilone.context[ecount][0] =
*(eptr - 1);
                epothilone.context[ecount][1] = *eptr;
                epothilone.context[ecount][2] = '-';
            } else {
                epothilone.context[ecount][0] =
*(eptr - 1);
                epothilone.context[ecount][1] = *eptr;
                epothilone.context[ecount][2] =
*(eptr + 1);
            }
        }
        epothilone.context[ecount][3] = '\0';
        eptr++;
        ecount++;
    }
    /* ----- align STARTER (current library entry) and

```

```

epothilone ----- */
    sptr = library[kk].monomersequence;
    lcount=0;
    while(*sptr != '\0') {
        fprintf(stdout,"library[%d].monomersequence[%d]=%c\n",
            kk,lcount,library[kk].monomersequence[lcount]);
        sptr++;
        lcount++;
    }
    /* Call maximal_adjacent_alignment until it no longer
returns more than two adjacent modules. There is no reason to
try to extract individual modules, because this is done as part of the recursive filling of spaces
from the library.
    */
    smallest_acceptable_piece = 2;
    eptr = epothilone.monomersequence;
    fprintf(stdout, "ALIGN_TARGET: ");
    while(*eptr != '\0') {
        fprintf(stdout,"%c",*eptr);
        eptr++;
    }
    fprintf(stdout,"\n");
    fprintf(stderr,"aligning %d %s\n",kk, library[kk].name);
    best_new_unmarked_entries_filled = 0;
    while((new_unmarked_entries_filled =
maximal_adjacent_alignment_and_dump(&epothilone,nwildcard,wildcards,library,kk,s
mallest_acceptable_piece)) >= STARTER_MINIMUM_ADJACENT_ALIGN) {
        if(best_new_unmarked_entries_filled < new_unmarked_entries_filled){
            best_new_unmarked_entries_filled =
new_unmarked_entries_filled;
        }
        if(DEBUG_STARTER) fprintf(stdout, "STARTER ALIGN:
new_unmarked_entries_filled=%d\n",new_unmarked_entries_filled);
        epothilonelen = strlen(epothilone.monomersequence);
        for(ii=0; ii< epothilonelen; ii++){

```



```

        if(DEBUG_STARTER) fprintf(stdout, "STARTER ALIGN:found
a best alignment
between epo.monomer[%d]=%c in library[%d].name=%s\n",
        ii,epothilone.monomersequence[ii],kk,epothilone.alignedPKSname[ii]);
    }
}
library[kk].recursion_tagged = TRUE;
fprintf(stdout,"ALIGN_TARGET:  \n");
dump_STARTER_align(epothilone,nwildcard,wildcards);
fprintf(stdout,"ALIGN_TARGET:  \n");
if(best_new_unmarked_entries_filled <= 1) {
    fprintf(stdout,"ALIGN_TARGET: PROBLEM
best_new_unmarked_entries_filled = %d\n",best_new_unmarked_entries_filled);
    fprintf(stdout,"ALIGN_TARGET: PROBLEM skipping this STARTER
entry for library[%d].name=%s\n",
        kk,library[kk].name);
    library[kk].recursion_tagged = FALSE;
    kk++;
    continue;
}

/* - fill in the gaps from the library - */
/* generate a fresh copy of epothilone in epotemp */
epothilonelen = strlen(epothilone.monomersequence);
nfilledmax = strlen(epothilone.monomersequence);
fprintf(stdout, "nfilledmax=%d\n", nfilledmax);
reset_epotemp(&epotemp,epothilone);
nfilled = 0;
for(ii=0; ii< epothilonelen; ii++) { if(epotemp.marked[ii] == TRUE) nfilled++; }
if(DEBUG_STARTER) fprintf(stdout,"nfilled from STARTER=%d\n",nfilled);
for(mmpass = 0; mmpass < nlib; mmpass++) {
    if(mmpass == kk) { continue; }
    reset_epotemp(&epotemp,epothilone);
    nfilled = 0;
    for(ii=0; ii< epothilonelen; ii++) {
if(epotemp.marked[ii] == TRUE) nfilled++; }

```

```

        if(nfilled >= nfilledmax) {
            output_fresh_alignment(&epotemp);
        } else {
            current_nmarked = nfilled;
            previous_nmarked = nfilled;
            smallest_acceptable_piece = 1;
            while((new_unmarked_entries_filled =
maximal_adjacent_alignment(&epotemp,nwildcard,wildcards,library,mmpass,smallest_
acceptable_piece)) >= MINIMUM_ADJACENT_ALIGN) {
                current_nmarked +=
new_unmarked_entries_filled;
                if(DEBUG_MATCH) fprintf(stdout,
"main: recursion_level=%d, mmpass=%d, previous_nmarked=%d,
current_nmarked=%d\n",
                recursion_counter,mmpass,previous_nmarked, current_nmarked);
            }
            nfilled = 0;
            for(ii=0; ii< epothilonelen; ii++) {
if(epotemp.marked[ii] == TRUE) nfilled++; }
            if(nfilled >= nfilledmax) {
                output_fresh_alignment(&epotemp);
                continue; /* no need to recurse */
            }
            if(DEBUG_MATCH) fprintf(stdout, "main:
about to RECURSE: mmpass=%d\n",mmpass);
            library[mmpass].recursion_tagged = TRUE;
            recursion_counter++;
            recurse_through_the_library(nfilledmax,epotemp,&epothilone,nwildcard,wildcards,
nlib,library,&recursion_counter);
            library[mmpass].recursion_tagged = FALSE;
            recursion_counter--;
        }
    }
    library[kk].recursion_tagged = FALSE;
    kk++;

```

```
    }/*nlib*/
}/*main*/

int  recurse_through_the_library(
int  nfilledmax,
LIB  epotemp,
LIB  *epothilone,
int  nwildcard,
char  wildcards[MAXWILD][MAXLEN],
int  nlib,
LIB  *library,
int  *recursion_counter)
{
int  ii=0;
int  ecoun=0,elen=0;
int  mmpass=0;
int  nfilled=0;
int  lcount=0;
int  previous_nmarked=0, current_nmarked=0;
int  smallest_acceptable_piece = 0;
char  *eptr;
char  *clibptr;
char  boundary[MAXNAMELEN];
int  new_unmarked_entries_filled=0;
LIB  epotemp_temp;
    if(DEBUG_MATCH) fprintf(stdout,"RECURSE: recursion_counter=%d,
nlib=%d\n",*recursion_counter,nlib);
    elen = strlen(epotemp.monomersequence);
    nfilled = 0;
    for(ii=0; ii< elen; ii++) { if(epotemp.marked[ii] == TRUE) nfilled++; }
    previous_nmarked = nfilled;
    current_nmarked = nfilled;
    smallest_acceptable_piece = 1;
    if(nfilled >= nfilledmax) { return 1; }
    for (mmpass = 0; mmpass < nlib; mmpass++) {
        if(*recursion_counter >= RECURSION_COUNTER_CUTOFF) {
```

```

        return 1;
    }
    if(DEBUG_MATCH) fprintf(stdout,"RECURSE:
recursion_counter=%d, mmpass=%d\n",*recursion_counter,mmpass);
    if(library[mmpass].recursion_tagged == TRUE) {
        if(DEBUG_MATCH) fprintf(stdout,"RECURSE:
library[%d].recursion_tagged=TRUE; skipping\n",mmpass);
        continue;
    }
    reset_epotemp(&epotemp_temp,epotemp);
    elen = strlen(epotemp_temp.monomersequence);
    nfilled = 0;
    for(ii=0; ii< elen; ii++) { if(epotemp_temp.marked[ii] ==
TRUE) nfilled++; }
    previous_nmarked = nfilled;
    current_nmarked = nfilled;
    while((new_unmarked_entries_filled =
maximal_adjacent_alignment(&epotemp_temp, nwildcard,wildcards,library,
mmpass,smallest_acceptable_piece)) >= 1) {
        current_nmarked += new_unmarked_entries_filled;
        if(DEBUG_MATCH) fprintf(stdout, "RECURSE: recursion_level=%d,
mmpass=%d, previous_nmarked=%d, current_nmarked=%d\n",
*recursion_counter,mmpass,previous_nmarked, current_nmarked);
    }
    elen = strlen(epotemp_temp.monomersequence);
    nfilled = 0;
    for(ii=0; ii< elen; ii++) { if(epotemp_temp.marked[ii] ==
TRUE) nfilled++; }
    if(nfilled >= nfilledmax) {
        output_fresh_alignment(&epotemp_temp);
        continue;
    }
    library[mmpass].recursion_tagged = TRUE;
    (*recursion_counter)++;
    recurse_through_the_library(nfilledmax,epotemp_temp,epothilone,nwildcard,wildca

```

```

rds,nlib,library,recursion_counter);
        library[mmpass].recursion_tagged = FALSE;
        (*recursion_counter)--;
    }/* mmpass */
}/*recurse_through_the_library*/
/*

    PURPOSE:
    INPUT:
    OUTPUT:
        returns the size of the largest maximal adjacent set of
monomers inserted.

    PROCEDURE:
*/
int maximal_adjacent_alignment(
LIB *epothilone,
int nwildcard,
char wildcards[MAXWILD][MAXLEN],
LIB *library,
int ilib,
int smallest_acceptable_piece)
{
int ii=0, jj=0, kk=0;
int ecount=0, lcount=0;
int epothilonelen=0;
int nlargestpiece=0, tnlargestpiece=0;
int hold_this_lcount=0, hold_this_ecount=0;
int wildcardmatch=FALSE;
char *wptr;
char *largestpiece_sptr, *largestpiece_eptr;
char *hold_this_place_eptr, *hold_this_place_sptr;
int largestpiece_lcount=0, largestpiece_ecount=0;
char *sptr, *eptr, *lptr, *bufptr;
if(DEBUG_WILDCARD) {
    if(nwildcard > 0) {
        fprintf(stdout, "maximal_:
```

```

wildcards[0]=%s\n",wildcards[0]);
    }
}

if(DEBUG_ALIGN) fprintf(stdout,"maximal_adjacent_alignment:
smallest_acceptable_piece=%d\n",smallest_acceptable_piece);

sptr = library[ilib].monomersequence;
eptr = epothilone->monomersequence;
ecount=0;
lcount=0;
nlargestpiece=0;
tnlargestpiece=0;
hold_this_place_eptr = eptr;
hold_this_ecount = ecount;
while (*eptr != '\0') {
    sptr = library[ilib].monomersequence;
    lcount = 0;
    hold_this_place_sptr = sptr;
    hold_this_lcount = lcount;
    wildcardmatch = FALSE;
    while(*sptr != '\0') {
        wildcardmatch = FALSE;
        if(epothilone->marked[ecount] == FALSE) {
            /* code for wildcards added MAS 05-16-00 */
            wptr = "";
            if(*eptr == 'X') { wptr = wildcards[0]; }
            else if(*eptr == 'Y') { wptr = wildcards[1]; }
            else if(*eptr == 'Z') { wptr = wildcards[2]; }
            while(*wptr != '\0') {
                if(*wptr == *sptr){
                    wildcardmatch = TRUE;
                    break;
                }
                wptr++;
            }
        }
        if((wildcardmatch == TRUE) || (*eptr ==

```

```
*sptr)) {  
    tnlargestpiece++;  
    if(DEBUG_ALIGN) fprintf(stdout,  
"FOUND a match: len=%d, epo(%d, %c), lib[%d].name=%s (%d, %c)\n",  
        tnlargestpiece,  
ecount, *eptr, ilib,library[ilib].name,lcount, *sptr);  
    if(tnlargestpiece > nlargestpiece) {  
        nlargestpiece = tnlargestpiece;  
        largestpiece_sptr =  
hold_this_place_sptr;  
  
        largestpiece_lcount =  
hold_this_lcount;  
  
        largestpiece_eptr =  
hold_this_place_eptr;  
  
        largestpiece_ecount =  
hold_this_ecount;  
  
        if(DEBUG_ALIGN)  
fprintf(stdout, "FOUND a largest piece: len=%d, epo(%d, %c),  
lib[%d].name=%s (%d, %c)\n",  
            nlargestpiece,  
largestpiece_ecount, *largestpiece_eptr,  
ilib,library[ilib].name,largestpiece_lcount, *largestpiece_sptr);  
    }  
    sptr++;  
    lcount++;  
    eptr++;  
    ecount++;  
} else {  
    tnlargestpiece = 0;  
    sptr++;  
    lcount++;  
    /* NEW */  
    hold_this_place_sptr = sptr;  
    hold_this_lcount = lcount;  
    eptr = hold_this_place_eptr;
```

```

        ecount = hold_this_ecount;
    }
    } else {
        tnlargestpiece = 0;
        break;
    }
}
tnlargestpiece = 0;
eptr = hold_this_place_eptr + 1;
ecount = hold_this_ecount + 1;
hold_this_place_eptr = eptr;
hold_this_ecount = ecount;
}

if(DEBUG_ALIGN) fprintf(stdout,"ALIGN: largest piece match is %d monomers
from %s\n",nlargestpiece,library[ilib].name);
if(DEBUG_ALIGN) fprintf(stdout,"ALIGN: largestpiece_ecount=%d,
largestpiece_lcount=%d\n",
    largestpiece_ecount, largestpiece_lcount);
if(nlargestpiece >= smallest_acceptable_piece) {
    if(DEBUG_ALIGN) fprintf(stdout,"ALIGN: incorporated\n");
    lcount = largestpiece_lcount;
    ecount = largestpiece_ecount;
    while(ecount < (nlargestpiece + largestpiece_ecount)) {
        epothilone->alignedsequence[ecount] =
library[ilib].monomersequence[lcount];
        strcpy(epothilone->alignedPKSname[ecount],library[ilib].name);
        strcpy(epothilone->context[ecount],library[ilib].context[lcount]);
        epothilone->marked[ecount] = TRUE;
        if(ecount < (nlargestpiece + largestpiece_ecount -
1)) epothilone->boundarytoright[ecount] = FALSE;
        lcount++;
        ecount++;
    }
}
}
return (nlargestpiece);

```



```
}/*maximal_adjacent_alignment*/
```

```
/*
```

```
    PURPOSE:
```

```
    INPUT:
```

```
    OUTPUT:
```

```
        returns the size of the largest maximal adjacent set of
monomers inserted.
```

```
    PROCEDURE:
```

```
*/
```

```
int    maximal_adjacent_alignment_and_dump(
```

```
LIB    *epothilone,
```

```
int    nwildcard,
```

```
char    wildcards[MAXWILD][MAXLEN],
```

```
LIB    *library,
```

```
int    ilib,
```

```
int    smallest_acceptable_piece)
```

```
{
```

```
int    ii=0, jj=0, kk=0;
```

```
int    ecount=0, lcount=0;
```

```
int    elen=0;
```

```
int    epothilonelen=0;
```

```
int    nlargestpiece=0, tnlargestpiece=0;
```

```
int    hold_this_lcount=0, hold_this_ecount=0;
```

```
int    wildcardmatch=FALSE;
```

```
char    *wptr;
```

```
char    *largestpiece_sptr, *largestpiece_eptr;
```

```
char    *hold_this_place_eptr, *hold_this_place_sptr;
```

```
int    largestpiece_lcount=0, largestpiece_ecount=0;
```

```
char    *sptr, *eptr, *lptr, *bufptr;
```

```
    if(DEBUG_WILDCARD) {
```

```
        if(nwildcard > 0) {
```

```
            fprintf(stdout, "maximal_:
```

```
wildcards[0]=%s\n", wildcards[0]);
```

```

    }
}

fprintf(stdout, "maximal_adjacent_alignment_and_dump:
smallest_acceptable_piece=%d\n", smallest_acceptable_piece);

sptr = library[ilib].monomersequence;
eptr = epothilone->monomersequence;
elen = strlen(epothilone->monomersequence);
ecount=0;
lcount=0;
nlargestpiece=0;
tnlargestpiece=0;
hold_this_place_eptr = eptr;
hold_this_ecount = ecount;
while (*eptr != '\0') {
    sptr = library[ilib].monomersequence;
    lcount = 0;
    /* NEW */
    hold_this_place_sptr = sptr;
    hold_this_lcount = lcount;
    wildcardmatch = FALSE;
    while(*sptr != '\0') {
        wildcardmatch = FALSE;
        if(epothilone->marked[ecount] == FALSE) {
            /* code for wildcards added MAS 05-16-00 */
            wptr = "";
            if(*eptr == 'X') { wptr = wildcards[0]; }
            else if(*eptr == 'Y') { wptr = wildcards[1]; }
            else if(*eptr == 'Z') { wptr = wildcards[2]; }
            while(*wptr != '\0') {
                if(*wptr == *sptr){
                    wildcardmatch = TRUE;
                    break;
                }
                wptr++;
            }
        }
    }
}

```

```

        if((wildcardmatch == TRUE) || (*eptr ==
*sptr)) {

            tnlargestpiece++;
            if(DEBUG_ALIGN) fprintf(stdout,
"FOUND a match: len=%d, epo(%d, %c), lib[%d].name=%s (%d, %c)\n",
                tnlargestpiece,
ecount, *eptr, ilib,library[ilib].name,lcount, *sptr);
            if(tnlargestpiece > nlargestpiece) {
                nlargestpiece = tnlargestpiece;
                largestpiece_sptr =
hold_this_place_sptr;

                largestpiece_lcount =
hold_this_lcount;

                largestpiece_eptr =
hold_this_place_eptr;

                largestpiece_ecount =
hold_this_ecount;

                if(DEBUG_ALIGN)
fprintf(stdout, "FOUND a largest piece: len=%d, epo(%d, %c), lib(%d, %c)\n",
                    nlargestpiece,
largestpiece_ecount, *largestpiece_eptr, largestpiece_lcount,
*largestpiece_sptr);
            }
            sptr++;
            lcount++;
            eptr++;
            ecount++;
        } else {
            tnlargestpiece = 0;
            sptr++;
            lcount++;
            hold_this_place_sptr = sptr;
            hold_this_lcount = lcount;
            /* NEW */
            eptr = hold_this_place_eptr;

```

```

        ecount = hold_this_ecount;
    }
} else {
    tnlargestpiece = 0;
    break;
}

}

tnlargestpiece = 0;
eptr = hold_this_place_eptr + 1;
ecount = hold_this_ecount + 1;
hold_this_place_eptr = eptr;
hold_this_ecount = ecount;
if(DEBUG_ALIGN) {
    fprintf(stdout, "incrementing
hold_this_place_eptr=%c,
hold_this_ecount=%d\n", *hold_this_place_eptr, hold_this_ecount);
}
}

if(DEBUG_ALIGN) fprintf(stdout, "ALIGN: largest piece match is %d monomers
from %s\n", nlargestpiece, library[ilib].name);
if(DEBUG_ALIGN) fprintf(stdout, "ALIGN: largestpiece_ecount=%d,
largestpiece_lcount=%d\n",
    largestpiece_ecount, largestpiece_lcount);
if(nlargestpiece >= smallest_acceptable_piece) {
    if(DEBUG_ALIGN) fprintf(stdout, "ALIGN: incorporated\n");
    lcount = largestpiece_lcount;
    ecount = largestpiece_ecount;
    fprintf(stdout, "ALIGN_TARGET: ");
    for(ii=0; ii<largestpiece_ecount; ii++) {
        fprintf(stdout, " ");
    }
    while(ecount < (nlargestpiece + largestpiece_ecount)) {
        epothilone->alignedsequence[ecount] =
library[ilib].monomersequence[lcount];
strcpy(epothilone->alignedPKSname[ecount], library[ilib].name);

```

```

        strcpy(epothilone->context[ecount],library[ilib].context[lcount]);
        epothilone->marked[ecount] = TRUE;
        if(ecount < (nlargestpiece + largestpiece_ecount -
1)) epothilone->boundarytoright[ecount] = FALSE;
        fprintf(stdout,"%c",library[ilib].monomersequence[lcount]);
        lcount++;
        ecount++;
    }
    for(ii=ecount; ii<elen; ii++) {
        fprintf(stdout," ");
    }
    fprintf(stdout,"%s\n",library[ilib].name);
}

return (nlargestpiece);
}/*maximal_adjacent_alignment_and_dump*/
int output_fresh_alignment(
LIB *epotemp)
{
    int ecount=0;
    char *eptr;
    char boundary[MAXNAMELEN];
    eptr = epotemp->monomersequence;
    ecount = 0;
    epotemp->nboundary = 0;
    strcpy(boundary,epotemp->alignedPKSname[ecount]);
    while(*eptr != '\0') {
        if(epotemp->boundarytoright[ecount] == TRUE) {
            epotemp->nboundary++;
        }
        ecount++;
        eptr++;
    }
    if(epotemp->nboundary > NBOUNDARY_CUTOFF) return 1;
    eptr = epotemp->monomersequence;
    ecount = 0;

```

```

fprintf(stdout,"HIT  ");
while(*eptr != '\0') {
    if(epotemp->alignedPKSname[ecount][0] == '\0') {
        if(epotemp->boundarytoright[ecount] == TRUE) {
            fprintf(stdout,"%c:TARG(%s)|",*eptr,epotemp->context[ecount]);
        } else {
            fprintf(stdout,"%c:TARG(%s)",*eptr,epotemp->context[ecount]);
        }
    } else {
        if(epotemp->boundarytoright[ecount] == TRUE) {
            fprintf(stdout,"%c:%4s(%s)|",*eptr,epotemp->alignedPKSname[ecount],epotemp->context[ecount]);
        } else {
            fprintf(stdout,"%c:%4s(%s)",*eptr,epotemp->alignedPKSname[ecount],epotemp->context[ecount]);
        }
    }
    ecount++;
    eptr++;
}
fprintf(stdout,"npiece %d",epotemp->nboundary);
fprintf(stdout,"\n");
return 1;
}

/*output_fresh_alignment*/
int  get_library(
char  *libraryfile,
LIB   *library)
{
int    ii=0,jj=0,kk=0,lcount=0;
int    nlib=0;
char   *bufptr,buf[MAXBUF];
char   tmonomersequence[MAXNAMELEN];
char   *lptr,*tptr;

```

```

FILE *libp;
for(kk=0; kk < MAX_LIB_ENTRIES; kk++) {
    library[kk].recursion_tagged = FALSE;
    for(ii=0; ii<MAXNAMELEN; ii++) {
        library[kk].name[ii] = '\0';
        library[kk].monomersequence[ii] = '\0';
        library[kk].alignedsequence[ii] = '\0';
        for(jj=0; jj<MAXLEN; jj++) {
            library[kk].alignedPKSname[jj][ii] = '\0';
            library[kk].marked[jj] = FALSE;
            library[kk].context[jj][0] = '\0';
            library[kk].context[jj][1] = '\0';
            library[kk].context[jj][2] = '\0';
            library[kk].context[jj][3] = '\0';
        }
    }
    library[kk].nboundary = 0;
}

/*    read in the library from PKS.lib */
if(NULL==(libp=fopen(libraryfile,"r"))) {
    fprintf(stdout,"TRY AGAIN; couldn't open %s\n",libraryfile);
    nlib=0;
    exit();
}
nlib=0;
while(nlib < MAX_LIB_ENTRIES) {
    if(NULL==fgets(buf,sizeof(buf),libp)) break;
    bufptr = buf;
    if(*bufptr == '#') continue;
    if(*bufptr != '\n') {
        lptr = library[nlib].name;
        while((*bufptr != ' ') && (*bufptr != '\0') &&
(*bufptr != '\n')){
            *lptr++ = *bufptr++;
        }
    }
}

```

```

        *lptr = '\0';
        if((*bufptr != '\0') && (*bufptr != '\n')) bufptr++;
        lptr = library[nlib].monomersequence;
        while((*bufptr != ' ') && (*bufptr != '\0') &&
(*bufptr != '\n')){
            /* This code specifically deletes inter-modular
double bonds, optional.*/

            if(*bufptr != '=') {
                *lptr++ = *bufptr++;
            } else {
                bufptr++;
            }
        }
        *lptr = '\0';
        if((*bufptr != '\0') && (*bufptr != '\n')) bufptr++;
        lptr = library[nlib].annotatedsequence;
        while((*bufptr != ' ') && (*bufptr != '\0') &&
(*bufptr != '\n')){
            *lptr++ = *bufptr++;
        }
        *lptr = '\0';
        if((*bufptr != '\0') && (*bufptr != '\n')) bufptr++;
        fprintf(stdout,"LIBRARY(%d) %s:
%s\n",nlib,library[nlib].name,library[nlib].monomersequence);
        fprintf(stdout,"LIBRARY(%d) %s:
%s\n",nlib,library[nlib].name,library[nlib].annotatedsequence);
        nlib++;
    }
}
fclose(libp);
for(kk=0; kk<nlib; kk++) {
    lcount = 0;
    lptr = library[kk].monomersequence;
    while(*lptr != '\0'){
        if(lcount == 0) {

```



```

        library[kk].context[lcount][0] = '-';
        library[kk].context[lcount][1] = *lptr;
        library[kk].context[lcount][2] = *(lptr + 1);
    } else {
        if(lcount ==
(strlen(library[kk].monomersequence) - 1)){
            library[kk].context[lcount][0] =
*(lptr - 1);
            library[kk].context[lcount][1] =
*lptr;
            library[kk].context[lcount][2] = '-';
        } else {
            library[kk].context[lcount][0] =
*(lptr - 1);
            library[kk].context[lcount][1] =
*lptr;
            library[kk].context[lcount][2] =
*(lptr + 1);
        }
    }
    library[kk].context[lcount][3] = '\0';
    lptr++;
    lcount++;
}
fprintf(stdout, "LIBRARY(%d) %s:
%s\n", kk, library[kk].name, library[kk].monomersequence);
fprintf(stdout, "LIBRARY(%d) %s:
%s\n", kk, library[kk].name, library[kk].annotatedsequence);
for(jj=0; jj< strlen(library[kk].monomersequence); jj++) {
    fprintf(stdout, "(%s)", library[kk].context[jj]);
}
fprintf(stdout, "\n");
}
return nlib;
}/*get_library*/

```

```
int    dump_STARTER_align(
LIB    epothilone,
int    nwildcard,
char    wildcards[MAXWILD][MAXLEN]
)
{
int    elen=0;
int    ecount=0,hold_ecount=0;
int    wildcardmatch=FALSE;
char    *sptr,*eptr,*wptr;
    elen = strlen(epothilone.monomersequence);
    eptr = epothilone.monomersequence;
    fprintf(stdout, "ALIGN_TARGET: ");
    while(*eptr != '\0') {
        fprintf(stdout,"%c",*eptr);
        eptr++;
    }
    fprintf(stdout,"\n");
    ecount=0;
    eptr = epothilone.monomersequence;
    sptr = epothilone.alignedsequence;
    fprintf(stdout, "ALIGN_TARGET: ");
    while(*eptr != '\0') {
        wildcardmatch = FALSE;
        wptr = "";
        if(*eptr == 'X') { wptr = wildcards[0]; }
        else if(*eptr == 'Y') { wptr = wildcards[1]; }
        else if(*eptr == 'Z') { wptr = wildcards[2]; }
        while(*wptr != '\0') {
            if(*wptr == *sptr){
                wildcardmatch = TRUE;
                break;
            }
            wptr++;
        }
    }
}
```

```

        if(wildcardmatch == TRUE) {
            fprintf(stdout, ".");
        } else {
            if(*eptr == *sptr) {
                fprintf(stdout, "|");
            } else {
                fprintf(stdout, " ");
            }
        }
        eptr++;
        sptr++;
    }
    fprintf(stdout, "\n");
    fprintf(stdout, "ALIGN_TARGET: ");
    eptr = epothilone.monomersequence;
    ecount=0;
    while(*eptr != '\0') {
        if(epothilone.alignedsequence[ecount] == '\0') {
            fprintf(stdout, " ");
        } else {
            fprintf(stdout, "%c", epothilone.alignedsequence[ecount]);
            hold_ecount = ecount;
        }
        eptr++;
        ecount++;
    }
    fprintf(stdout, "
%s\n", epothilone.alignedPKSname[hold_ecount]);
    fprintf(stdout, "STARTER_ALIGN:\n");
} /*dump_STARTER_align*/
reset_epotemp(
LIB    *epotemp,
LIB    epothilone)
{
    int    jj=0, elen=0;

```

```

elen = strlen(epothilone.monomersequence);
strcpy(epotemp->name,epothilone.name);
strcpy(epotemp->monomersequence,epothilone.monomersequence);
strcpy(epotemp->alignedsequence,epothilone.alignedsequence);
epotemp->nboundary = 0;
for(jj=0; jj< elen; jj++) {
strcpy(epotemp->alignedPKSname[jj],epothilone.alignedPKSname[jj]);
    epotemp->marked[jj] = epothilone.marked[jj];
    strcpy(epotemp->context[jj],epothilone.context[jj]);
    epotemp->boundarytoright[jj] = epothilone.boundarytoright[jj];
}
}/*reset_epotemp*/

```

EXAMPLE 6

Source Code:

```

#include <stdio.h>

/*      ~siani/programs/morph/morph4.c
PURPOSE:  To recursively traverse all the entries in PKS.lib, generating all feasible
combinations of PKS modules to make the TARGET (e.g., epothilone).
INPUT:    -b number_boundary_cutoff: lets user set the maximum number of
boundaries in output lines. This defaults to 5 (#define NBOUNDARY_CUTOFF 5) which is a
reasonable assumption for something of the length of epothilone (8 modules). However, when
looking at disco-dermolide which has 11 modules, a cutoff of 5 sometimes results in too few
output lines; it is too restrictive.
-d allows one to ignore the inter-modular doublebonds in the library file.
-l libraryfile: tab-delimited CHUCKLES-coded polyketides file with the following
columns
    1. polyketide name
    2. plain CHUCKLES
    3. annotated CHUCKLES (contains information about post-synthetic
modifications)
    4. source organism
-n targetname: user-defined name (e.g., epoD)
-t targetsequence: CHUCKLES-coded polyketide of desired TARGET (e.g.,
MEMLJGE)
-w, -x, -y, -z sets of wildcards: sets of monomers for particular positions appearing in
targetsequence. The wildcards can effectively be used for analoging the TARGET polyketide.

```

Hard-coded parameters which may be reset (requires recompiling):

```
#define NBOUNDARY_CUTOFF 5
```

NBOUNDARY_CUTOFF determines the maximum number of non-native inter-

modular interfaces which are contained in the output. This is now set to 5, but may be increased when the user does not care about inefficiencies introduced by these interfaces or when the targetsequence is very lengthy.

```
#define RECURSION_COUNTER_CUTOFF      2
```

RECURSION_COUNTER_CUTOFF specifies the number of levels of recursion (defaults to 0, 1, 2) acceptable for the run. This limit must be set since the large PKS library can result in recursion that will combinatorially explode. Because of the multi-directionality of the alignments (using every library entry as a STARTER), there is no need to go beyond 2 levels of recursion. However, there may be cases in the future where this number should be increased. Note that while recursion will eventually terminate without this parameter, runs with a library over about 20 PKS entries may run for years on a reasonably fast computer.

OUTPUT: All combinations of modules that meet parameters set by user.

Example output from MEMLJDGE (epothilone D) using subset of PKS.lib.

Vertical bars indicate non-native inter-modular interfaces. Last column contains the number of "pieces" that are needed to put together the PKS.

Names of PKSs have been abbreviated to fit them in these comments.

```
HIT M:3atyl(FMN)| E:tedan(GEH)| M:aldga(BML) L:aldga(MLG)| J:aldga(GJD)
D:aldga(JDL)| G:tedan(JGE) E:tedan(GEH)| 5
HIT M:albM1(LME) E:albM1(MEJ)| M:albM1(LME)| L:aldga(MLG)| J:aldga(GJD)
D:aldga(JDL)| G:tedan(JGE) E:tedan(GEH)| 5
HIT M:albM1(LME) E:albM1(MEJ)| M:aldga(BML) L:aldga(MLG)| J:aldga(GJD)
D:aldga(JDL)| G:3atyl(NGO)| E:albM1(MEJ)| 5
HIT M:albM1(LME) E:albM1(MEJ)| M:aldga(BML) L:aldga(MLG)| J:aldga(GJD)
D:aldga(JDL)| G:aldga(LGJ)| E:albM1(MEJ)| 5
HIT M:albM1(LME) E:albM1(MEJ)| M:aldga(BML) L:aldga(MLG)| J:aldga(GJD)
D:aldga(JDL)| G:aldga(LGJ)| E:albM1(MEJ)| 5
```

USAGE:

```
morph3 -l libraryfile -n targetname -t targetsequence [-w W-wildcards] [-x X-
wildcards] [-y Y-wildcards] [-z Z-wildcards] -d
```

examples:

```
# generate combinations that yield epothilone D
```

```
%morph3 -l PKS.lib -n epoD -t MEMLJDGE > omorph3_epoD
```

```
%grep HIT omorph3_epoD | sort | uniq | sort +10 -11 >
```

```

omorph3_epoD.uniq.sort
    %egrep ALIGN_TARGET omorph3_epoD >
omorph3_epoD_STARTER_ALIGN
    # generate combinations that yield epothilone D    with a C13-hydroxyl
    %morph3 -l PKS.lib -n epoD-13OH -t MEXLJDGE -x ABCD > oepoD-13OH
    %egrep HIT oepoD-13OH | sort | uniq | sort +10 -11 > oepoD-13OH.uniq.sort
    %egrep ALIGN_TARGET oepoD-13OH > oepoD-13OH_STARTER_ALIGN
    # generate combination that yield epothilone with the following wildcards (set 1)
    %morph3 -l PKS.lib -n epoD-set1 -t MEXYZDgE -x ABCD -y LBFIN -z
JACGM > oepoD-set1
    %grep HIT oepoD-set1 | sort | uniq | sort +10 -11 > oepoD-set1.uniq.sort

    # generate combination that yield epothilone with the following wildcards (set 2)
    %morph3 -l PKS.lib -n epoD-set2 -t MEXYZDgE -x JK -y EF -z JACGM >
oepoD-set2
    %grep HIT oepoD-set2 | sort | uniq | sort +10 -11 > oepoD-set2.uniq.sort

```

LIMITATIONS:

Current implementation cannot handle intra-modular modifications/splitting because morph is operating at the monomer level. Future implementations could convert the CHUCKLES-encoded strings into the corresponding and equivalent SMILES and then perform more complex chemical analysis of the PKS molecular graphs. Currently, inter-modular double bonds are present in the library, but are ignored by the morph program.

MODIFICATIONS:

- + added ability to include user-defined wildcards (X, Y, or Z) on the command line. MAS 05-16-00.
- + added additional wildcard (W). MAS 05-30-00.
- + added addition (summary) column to HIT output list. MAS 05-30-00.
- + added command line argument for suppressing the inter-modular double bonds in the library. Default is not to treat these as separate modules. MAS 05-31-00.
- + added column that contains the length of the largest matching fragment MAS

06-05-00

*/

```
#define TRUE      1
#define FALSE     0
#define DEBUG_MATCH FALSE
#define DEBUG_STARTER FALSE
#define DEBUG_ALIGN FALSE
#define DEBUG_RECURSE FALSE
#define DEBUG_WILDCARD FALSE

#define MAXLEN      80
#define MAX_TYL_LEN  6
#define MAX_EPO_LEN  6
#define MAXNAMELEN  500
#define MAX_LIB_ENTRIES 500
#define MAXWILD  4
#define MAXBUF  1000

#define NBOUNDARY_CUTOFF 5
#define RECURSION_COUNTER_CUTOFF 2
#define STARTER_MINIMUM_ADJACENT_ALIGN 2
#define MINIMUM_ADJACENT_ALIGN 2

typedef struct _lib {
    char  name[MAXNAMELEN];
    char  monomersequence[MAXNAMELEN];
    char  annotatedsequence[MAXNAMELEN];
    char  alignedsequence[MAXNAMELEN];
    char  alignedPKSname[MAXLEN][MAXNAMELEN];
    int   boundarytoright[MAXNAMELEN];
    int   marked[MAXLEN];
    char  context[MAXLEN][4];
    int   recursion_tagged;
    int   nboundary;
} LIB;
```

```
main(int argc, char **argv)
{
    int    ii=0, jj=0, kk=0, ll=0;
    int    nlib=0;
    int    ecoun=0;
    int    nfilled=0, nfilledmax=0;
    int    epothilonelen=0;
    int    nlargestpiece=0, tnlargestpiece=0;
    int    mmpass=0;
    int    lcount=0;
    int    new_unmarked_entries_filled=0;
    int    recursion_counter = 0;
    int    nwildcard=0;
    int    best_new_unmarked_entries_filled = 0;
    int    smallest_acceptable_piece = 0;
    int    current_nmarked=0, previous_nmarked=0;
    int    inter_modular_db_flag_off=FALSE;

    int    nboundary_cutoff=NBOUNDARY_CUTOFF;

    char    *sptr, *eptr, *lptr, *bufptr;
    char    *clibptr;
    char    *libraryfile;
    char    *targetsequence, *targetname;
    char    buf[MAXBUF];
    char    wildcards[MAXWILD][MAXLEN];

    FILE    *libp;

    LIB    epotemp;
    LIB    library[MAX_LIB_ENTRIES];
    LIB    epothilone;

    char    *programe;
    char    **filelist, **fileptr;
```



```

libraryfile = "";
targetsequence = "";
targetname = "";

for(ii=0; ii<MAXWILD; ii++) {
    for(jj=0; jj<MAXLEN; jj++) {
        wildcards[ii][jj] = '\0';
    }
}

/* process arguments */
filelist = fileptr = (char **)(malloc(argc * sizeof(*argv)));

programe = *argv++;
if(argc<2) {
    fprintf(stderr,"usage:%s [-b nboundary_cutoff] [-d] -l libraryfile -n targetname -t
targetsequence [-w W-wildcards] [-x X-wildcards] [-y Y-wildcards] [-z Z-wildcards]
\n",programe);
    exit();
}

while(argc-- > 1) {
    if(argv[0][0] == '-' && argv[0][1] != '\0') {
        /* handle option */
        *++(*argv); /* advance past the minus */
        switch(**argv) {
            case 'b': /* get number of boundaries cutoff for output of
alignments */
                argv++; argc--;
                sscanf(argv[0],"%d",&nboundary_cutoff);
                fprintf(stderr,"-b:
nboundary_cutoff=%d\n",nboundary_cutoff);
                break;
            case 'd': /* ignore inter-modular double bonds in the library file */
                inter_modular_db_flag_off = TRUE;

```

```
    fprintf(stderr, "-d: inter-modular double bonds ignored.\n");
    break;
case 'l': /* get library input filename (PKS.lib) */
    argv++; argc--;
    libraryfile = argv[0];
    fprintf(stderr, "-l: libraryfile=%s\n", libraryfile);
    break;
case 'n': /* get target name string */
    argv++; argc--;
    targetname = argv[0];
    fprintf(stderr, "-t: targetname=%s\n", targetname);
    break;
case 't': /* get target sequence string */
    argv++; argc--;
    targetsequence = argv[0];
    fprintf(stderr, "-t: targetsequence=%s\n", targetsequence);
    break;
case 'w': /* get a wildcard string */
    argv++; argc--;
    strcpy(wildcards[0], argv[0]);
    fprintf(stderr, "-w: wildcards[%d]=%s\n", 0, wildcards[0]);
    nwildcard++;
    break;
case 'x': /* get a wildcard string */
    argv++; argc--;
    strcpy(wildcards[1], argv[0]);
    fprintf(stderr, "-x: wildcards[%d]=%s\n", 1, wildcards[1]);
    nwildcard++;
    break;
case 'y': /* get a wildcard string */
    argv++; argc--;
    strcpy(wildcards[2], argv[0]);
    fprintf(stderr, "-y: wildcards[%d]=%s\n", 2, wildcards[2]);
    nwildcard++;
    break;
```

```
case 'z': /* get a wildcard string */
    argv++; argc--;
    strcpy(wildcards[3],argv[0]);
    fprintf(stderr,"-z: wildcards[%d]=%s\n",3,wildcards[3]);
    nwildcard++;
    break;

case 'W': /* get a wildcard string */
    argv++; argc--;
    strcpy(wildcards[0],argv[0]);
    fprintf(stderr,"-w: wildcards[%d]=%s\n",0,wildcards[0]);
    nwildcard++;
    break;

case 'X': /* get a wildcard string */
    argv++; argc--;
    strcpy(wildcards[1],argv[0]);
    fprintf(stderr,"-x: wildcards[%d]=%s\n",1,wildcards[1]);
    nwildcard++;
    break;

case 'Y': /* get a wildcard string */
    argv++; argc--;
    strcpy(wildcards[2],argv[0]);
    fprintf(stderr,"-y: wildcards[%d]=%s\n",2,wildcards[2]);
    nwildcard++;
    break;

case 'Z': /* get a wildcard string */
    argv++; argc--;
    strcpy(wildcards[3],argv[0]);
    fprintf(stderr,"-z: wildcards[%d]=%s\n",3,wildcards[3]);
    nwildcard++;
    break;

default:
    fprintf(stderr,"%s unknown option; ignored\n",*argv);
} /*switch*/

} else { /* a regular filename */
    *fileptr++ = *argv;
```

```
        *fileptr = NULL;
    }

    argv++;
}/*while*/

if(nwildcard > 0) {
    for(ii=0; ii<nwildcard; ii++) {
        fprintf(stderr,"wildcards[%d]=%s\n",ii,wildcards[ii]);
        fprintf(stdout,"wildcards[%d]=%s\n",ii,wildcards[ii]);
    }
}

epothilone.nboundary = 0;
for(ii=0; ii<MAXNAMELEN; ii++) {
    epothilone.name[ii] = '\0';
    epothilone.monomersequence[ii] = '\0';
    epothilone.alignedsequence[ii] = '\0';
    epothilone.boundarytoright[ii] = TRUE;

    for(jj=0; jj<MAXLEN; jj++) {
        epothilone.alignedPKSname[jj][ii] = '\0';
        epothilone.marked[jj] = FALSE;
        epothilone.context[jj][0] = '\0';
        epothilone.context[jj][1] = '\0';
        epothilone.context[jj][2] = '\0';
    }
}

strcpy(epothilone.name,targetname);
strcpy(epothilone.monomersequence,targetsequence);

fprintf(stdout, "TARGET:   %s\n", epothilone.monomersequence);

ecount = 0;
```

```
eptr = epothilone.monomersequence;
while(*eptr != '\0'){
    if(ecount == 0) {
        epothilone.context[ecount][0] = '-';
        epothilone.context[ecount][1] = *eptr;
        epothilone.context[ecount][2] = *(eptr + 1);
    } else {
        if(ecount == (strlen(epothilone.monomersequence) - 1)){
            epothilone.context[ecount][0] = *(eptr - 1);
            epothilone.context[ecount][1] = *eptr;
            epothilone.context[ecount][2] = '-';
        } else {
            epothilone.context[ecount][0] = *(eptr - 1);
            epothilone.context[ecount][1] = *eptr;
            epothilone.context[ecount][2] = *(eptr + 1);
        }
    }
    epothilone.context[ecount][3] = '\0';
    eptr++;
    ecount++;
}

for(ii=0; ii<ecount; ii++) {
    fprintf(stdout,"%s\n",epothilone.context[ii]);
}

/* ----- library ----- */
nlib = get_library(libraryfile,library,inter_modular_db_flag_off);
fprintf(stdout,"nlib=%d\n",nlib);

kk=0; while(kk < nlib) {
```

```

/* ----- zero out the epothilone entry with respect to a new alignment ----- */
for(ii=0; ii<MAXNAMELEN; ii++) {
    epothilone.alignedsequence[ii] = '\0';
    epothilone.boundarytoright[ii] = TRUE;
    for(jj=0; jj<MAXLEN; jj++) {
        epothilone.alignedPKSname[jj][ii] = '\0';
        epothilone.marked[jj] = FALSE;
    }
}

/* ----- reset the context back to that in epothilone ----- */
ecount = 0;
eptr = epothilone.monomersequence;
while(*eptr != '\0'){
    if(ecount == 0) {
        epothilone.context[ecount][0] = '-';
        epothilone.context[ecount][1] = *eptr;
        epothilone.context[ecount][2] = *(eptr + 1);
    } else {
        if(ecount == (strlen(epothilone.monomersequence) - 1)){
            epothilone.context[ecount][0] = *(eptr - 1);
            epothilone.context[ecount][1] = *eptr;
            epothilone.context[ecount][2] = '-';
        } else {
            epothilone.context[ecount][0] = *(eptr - 1);
            epothilone.context[ecount][1] = *eptr;
            epothilone.context[ecount][2] = *(eptr + 1);
        }
    }
    epothilone.context[ecount][3] = '\0';
    eptr++;
    ecount++;
}

```

```

/* ----- align STARTER (current library entry) and epothilone ----- */

```

```

sptr = library[kk].monomersequence;
lcount=0;
while(*sptr != '\0') {
    fprintf(stdout,"library[%d].monomersequence[%d]=%c\n",
            kk,lcount,library[kk].monomersequence[lcount]);
    sptr++;
    lcount++;
}

/* Call maximal_adjacent_alignment until it no longer returns more than
   two adjacent modules. There is really no reason to try to extract
   individual modules because this will be done as part of the
   recursive filling of spaces from the library.
*/
smallest_acceptable_piece = 2;
eptr = epothilone.monomersequence;
fprintf(stdout, "ALIGN_TARGET: ");
while(*eptr != '\0') {
    fprintf(stdout,"%c",*eptr);
    eptr++;
}
fprintf(stdout,"\n");
fprintf(stderr,"aligning %d %s\n",kk, library[kk].name);

best_new_unmarked_entries_filled = 0;
while((new_unmarked_entries_filled =
maximal_adjacent_alignment_and_dump(&epothilone,nwildcard,wildcards,library,kk,smallest_
acceptable_piece)) >= STARTER_MINIMUM_ADJACENT_ALIGN) {

    if(best_new_unmarked_entries_filled < new_unmarked_entries_filled){
        best_new_unmarked_entries_filled =
new_unmarked_entries_filled;
    }
    if(DEBUG_STARTER) fprintf(stdout, "STARTER ALIGN:

```

```
new_unmarked_entries_filled=%d\n",new_unmarked_entries_filled);

    epothilonelen = strlen(epothilone.monomersequence);
    for(ii=0; ii< epothilonelen; ii++){
        if(DEBUG_STARTER) fprintf(stdout, "STARTER ALIGN:found
a best alignment between epo.monomer[%d]=%c in library[%d].name=%s\n",

            ii,epothilone.monomersequence[ii],kk,epothilone.alignedPKSname[ii]);
    }
}

library[kk].recursion_tagged = TRUE;

fprintf(stdout,"ALIGN_TARGET:  \n");
dump_STARTER_align(epothilone,nwildcard,wildcards);
fprintf(stdout,"ALIGN_TARGET:  \n");

if(best_new_unmarked_entries_filled <= 1) {
    fprintf(stdout,"ALIGN_TARGET: PROBLEM
best_new_unmarked_entries_filled = %d\n",best_new_unmarked_entries_filled);
    fprintf(stdout,"ALIGN_TARGET: PROBLEM skipping this STARTER
entry for library[%d].name=%s\n",
        kk,library[kk].name);
    library[kk].recursion_tagged = FALSE;
    kk++;
    continue;
}

/* ----- fill in the gaps from the library ----- */

/* generate a fresh copy of epothilone in epotemp */
epothilonelen = strlen(epothilone.monomersequence);
nfilledmax = strlen(epothilone.monomersequence);
```



```

fprintf(stdout, "nfilledmax=%d\n", nfilledmax);

reset_epotemp(&epotemp, epothilone);
nfilled = 0;
for(ii=0; ii< epothilonelen; ii++) { if(epotemp.marked[ii] == TRUE) nfilled++; }
if(DEBUG_STARTER) fprintf(stdout, "nfilled from STARTER=%d\n", nfilled);

for(mmpass = 0; mmpass < nlib; mmpass++) {

    if(mmpass == kk) { continue; }
    reset_epotemp(&epotemp, epothilone);

    nfilled = 0;
    for(ii=0; ii< epothilonelen; ii++) { if(epotemp.marked[ii] == TRUE)
nfilled++; }

    if(nfilled >= nfilledmax) {
        output_fresh_alignment(&epotemp, nboundary_cutoff);
    } else {

        current_nmarked = nfilled;
        previous_nmarked = nfilled;
        smallest_acceptable_piece = 2;
        while((new_unmarked_entries_filled =
maximal_adjacent_alignment(&epotemp, nwildcard, wildcards, library, mmpass, smallest_acceptab
le_piece)) >= MINIMUM_ADJACENT_ALIGN) {
            current_nmarked += new_unmarked_entries_filled;
            if(DEBUG_MATCH) fprintf(stdout, "main:
recursion_level=%d, mmpass=%d, previous_nmarked=%d, current_nmarked=%d\n",
recursion_counter, mmpass, previous_nmarked,
current_nmarked);
        }
    }
}

```

```

        nfilled = 0;
        for(ii=0; ii< epothilonelen; ii++) { if(epotemp.marked[ii] ==
TRUE) nfilled++; }

        if(nfilled >= nfilledmax) {
            output_fresh_alignment(&epotemp,nboundary_cutoff);
            continue; /* no need to recurse */
        }

        if(DEBUG_MATCH) fprintf(stdout, "main: about to RECURSE:
mmpass=%d\n",mmpass);

        library[mmpass].recursion_tagged = TRUE;
        recursion_counter++;

        recurse_through_the_library(nfilledmax,epotemp,&epothilone,nwildcard,wildcards,nlib,l
ibrary,&recursion_counter,nboundary_cutoff);

        library[mmpass].recursion_tagged = FALSE;
        recursion_counter--;
    }
}

    library[kk].recursion_tagged = FALSE;
    kk++;
}/*nlib*/

}/*main*/

/*
-----
PURPOSE:
INPUT:
OUTPUT:
PROCEDURE:
*/
int  recurse_through_the_library(

```

```
int    nfilledmax,
LIB    epotemp,
LIB    *epothilone,
int    nwildcard,
char   wildcards[MAXWILD][MAXLEN],
int    nlib,
LIB    *library,
int    *recursion_counter,
int    nboundary_cutoff)
{
int    ii=0;
int    ecount=0,elen=0;
int    mmpass=0;
int    nfilled=0;
int    lcount=0;
int    previous_nmarked=0, current_nmarked=0;
int    smallest_acceptable_piece = 0;
char   *eptr;
char   *clibptr;
char   boundary[MAXNAMELEN];
int    new_unmarked_entries_filled=0;
LIB    epotemp_temp;

    if(DEBUG_MATCH) fprintf(stdout,"RECURSE: recursion_counter=%d,
nlib=%d\n",*recursion_counter,nlib);

    elen = strlen(epotemp.monomersequence);
    nfilled = 0;
    for(ii=0; ii< elen; ii++) { if(epotemp.marked[ii] == TRUE) nfilled++; }
    previous_nmarked = nfilled;
    current_nmarked = nfilled;
    smallest_acceptable_piece = 1;
    if(nfilled >= nfilledmax) { return 1; }

    for (mmpass = 0; mmpass < nlib; mmpass++) {
```

```
        if(*recursion_counter >= RECURSION_COUNTER_CUTOFF) {
            return 1;
        }

        if(DEBUG_MATCH) fprintf(stdout, "RECURSE: recursion_counter=%d,
mmpass=%d\n", *recursion_counter, mmpass);
        if(library[mmpass].recursion_tagged == TRUE) {
            if(DEBUG_MATCH) fprintf(stdout, "RECURSE:
library[%d].recursion_tagged=TRUE; skipping\n", mmpass);
            continue;
        }

        reset_epotemp(&epotemp_temp, epotemp);

        elen = strlen(epotemp_temp.monomersequence);
        nfilled = 0;
        for(ii=0; ii< elen; ii++) { if(epotemp_temp.marked[ii] == TRUE) nfilled++; }
        previous_nmarked = nfilled;
        current_nmarked = nfilled;

        while((new_unmarked_entries_filled =
maximal_adjacent_alignment(&epotemp_temp, nwildcard, wildcards, library,
mmpass, smallest_acceptable_piece)) >= 1) {
            current_nmarked += new_unmarked_entries_filled;
            if(DEBUG_MATCH) fprintf(stdout, "RECURSE: recursion_level=%d,
mmpass=%d, previous_nmarked=%d, current_nmarked=%d\n",
                *recursion_counter, mmpass, previous_nmarked,
current_nmarked);
        }

        elen = strlen(epotemp_temp.monomersequence);
        nfilled = 0;
        for(ii=0; ii< elen; ii++) { if(epotemp_temp.marked[ii] == TRUE) nfilled++; }
```

```

        if(nfilled >= nfilledmax) {
            output_fresh_alignment(&epotemp_temp,nboundary_cutoff);
            continue;
        }

        library[mmpass].recursion_tagged = TRUE;
        (*recursion_counter)++;

        recurse_through_the_library(nfilledmax,epotemp_temp,epothilone,nwildcard,wildcards,
nlib,library,recursion_counter,nboundary_cutoff);
        library[mmpass].recursion_tagged = FALSE;
        (*recursion_counter)--;

    }/* mmpass */

}/*recurse_through_the_library*/

/* -----
PURPOSE:
INPUT:
OUTPUT:
        returns the size of the largest maximal adjacent set of monomers inserted.
PROCEDURE:
*/
int maximal_adjacent_alignment(
LIB *epothilone,
int nwildcard,
char wildcards[MAXWILD][MAXLEN],
LIB *library,
int ilib,
int smallest_acceptable_piece)
{
    int ii=0, jj=0, kk=0;
    int ecount=0, lcount=0;

```

```
int    epothilonelen=0;
int    nlargestpiece=0,tnlargestpiece=0;
int    hold_this_lcount=0, hold_this_ecount=0;
int    wildcardmatch=FALSE;
char    *wptr;
char    *largestpiece_sptr,*largestpiece_eptr;
char    *hold_this_place_eptr, *hold_this_place_sptr;
int    largestpiece_lcount=0, largestpiece_ecount=0;

char    *sptr, *eptr, *lptr,*bufptr;

    if(DEBUG_WILDCARD) {
        if(nwildcard > 0) {
            fprintf(stdout,"maximal_ : wildcards[0]=%s\n",wildcards[0]);
        }
    }

    if(DEBUG_ALIGN) fprintf(stdout,"maximal_adjacent_alignment:
smallest_acceptable_piece=%d\n",smallest_acceptable_piece);

    sptr = library[ilib].monomersequence;
    eptr = epothilone->monomersequence;
    ecount=0;
    lcount=0;
    nlargestpiece=0;
    tnlargestpiece=0;
    hold_this_place_eptr = eptr;
    hold_this_ecount = ecount;
    while (*eptr != '\0') {

        sptr = library[ilib].monomersequence;
        lcount = 0;
        hold_this_place_sptr = sptr;
        hold_this_lcount = lcount;
        wildcardmatch = FALSE;
```

```

while(*sptr != '\0') {
    wildcardmatch = FALSE;
    if(epothilone->marked[ecount] == FALSE) {
        /* code for wildcards added MAS 05-16-00 */
        wptr = "";
        if(*eptr == 'W') { wptr = wildcards[0]; }
        else if(*eptr == 'X') { wptr = wildcards[1]; }
        else if(*eptr == 'Y') { wptr = wildcards[2]; }
        else if(*eptr == 'Z') { wptr = wildcards[3]; }

        while(*wptr != '\0') {
            if(*wptr == *sptr){
                wildcardmatch = TRUE;
                break;
            }
            wptr++;
        }
        if((wildcardmatch == TRUE) || (*eptr == *sptr)) {

            tnlargestpiece++;
            if(DEBUG_ALIGN) fprintf(stdout, "FOUND a match:
len=%d, epo(%d, %c), lib[%d].name=%s (%d, %c)\n",
                                tnlargestpiece, ecount, *eptr,
                                ilib, library[ilib].name, lcount, *sptr);

            if(tnlargestpiece > nlargestpiece) {
                nlargestpiece = tnlargestpiece;
                largestpiece_sptr = hold_this_place_sptr;
                largestpiece_lcount = hold_this_lcount;
                largestpiece_eptr = hold_this_place_eptr;
                largestpiece_ecount = hold_this_ecount;
                if(DEBUG_ALIGN) fprintf(stdout, "FOUND a
largest piece: len=%d, epo(%d, %c), lib[%d].name=%s (%d, %c)\n",
                                nlargestpiece, largestpiece_ecount,
                                *largestpiece_eptr, ilib, library[ilib].name, largestpiece_lcount, *largestpiece_sptr);
            }

```

```
        sptr++;
        lcount++;
        eptr++;
        ecount++;
    } else {
        tnlargestpiece = 0;
        sptr++;
        lcount++;
        /* NEW */
        hold_this_place_sptr = sptr;
        hold_this_lcount = lcount;
        eptr = hold_this_place_eptr;
        ecount = hold_this_ecount;
    }
} else {
    tnlargestpiece = 0;
    break;
}
}

tnlargestpiece = 0;
eptr = hold_this_place_eptr + 1;
ecount = hold_this_ecount + 1;
hold_this_place_eptr = eptr;
hold_this_ecount = ecount;
}

if(DEBUG_ALIGN) fprintf(stdout,"ALIGN: largest piece match is %d monomers from
%s\n",nlargestpiece,library[ilib].name);

if(DEBUG_ALIGN) fprintf(stdout,"ALIGN: largestpiece_ecount=%d,
largestpiece_lcount=%d\n",
    largestpiece_ecount,largestpiece_lcount);

if(nlargestpiece >= smallest_acceptable_piece) {
    if(DEBUG_ALIGN) fprintf(stdout,"ALIGN: incorporated\n");
    lcount = largestpiece_lcount;
    ecount = largestpiece_ecount;
```



```

        while(ecount < (nlargestpiece + largestpiece_ecount)) {

            epothilone->alignedsequence[ecount] =
library[ilib].monomersequence[lcount];
            strcpy(epothilone->alignedPKSname[ecount],library[ilib].name);
            strcpy(epothilone->context[ecount],library[ilib].context[lcount]);
            epothilone->marked[ecount] = TRUE;
            if(ecount < (nlargestpiece + largestpiece_ecount - 1)) epothilone-
>boundarytoright[ecount] = FALSE;

            lcount++;
            ecount++;

        }
    }

```

```

    return (nlargestpiece);

```

```

}/*maximal_adjacent_alignment*/

```

```

/* -----
PURPOSE:
INPUT:
OUTPUT:
    returns the size of the largest maximal adjacent set of monomers inserted.
PROCEDURE:
*/
int  maximal_adjacent_alignment_and_dump(
LIB  *epothilone,
int  nwildcard,
char wildcards[MAXWILD][MAXLEN],
LIB  *library,
int  ilib,
int  smallest_acceptable_piece)
{

```

```
int    ii=0, jj=0, kk=0;
int    ecoun=0, lcount=0;
int    elen=0;
int    epothilonelen=0;
int    nlargestpiece=0, tnlargestpiece=0;
int    hold_this_lcount=0, hold_this_ecoun=0;
int    wildcardmatch=FALSE;
char    *wptr;
char    *largestpiece_sptr, *largestpiece_eptr;
char    *hold_this_place_eptr, *hold_this_place_sptr;
int    largestpiece_lcount=0, largestpiece_ecoun=0;

char    *sptr, *eptr, *lptr, *bufptr;

    if(DEBUG_WILDCARD) {
        if(nwildcard > 0) {
            fprintf(stdout, "maximal_: wildcards[0]=%s\n", wildcards[0]);
        }
    }

    fprintf(stdout, "maximal_adjacent_alignment_and_dump:
smallest_acceptable_piece=%d\n", smallest_acceptable_piece);

    sptr = library[ilib].monomersequence;
    eptr = epothilone->monomersequence;
    elen = strlen(epothilone->monomersequence);
    ecoun=0;
    lcount=0;
    nlargestpiece=0;
    tnlargestpiece=0;
    hold_this_place_eptr = eptr;
    hold_this_ecoun = ecoun;
    while (*eptr != '\0') {

        sptr = library[ilib].monomersequence;
        lcount = 0;
```

```

/* NEW */
hold_this_place_sptr = sptr;
hold_this_lcount = lcount;
wildcardmatch = FALSE;

while(*sptr != '\0') {
    wildcardmatch = FALSE;
    if(epothilone->marked[ecount] == FALSE) {
        /* code for wildcards added MAS 05-16-00 */
        wptr = "";
        if(*eptr == 'W') { wptr = wildcards[0]; }
        else if(*eptr == 'X') { wptr = wildcards[1]; }
        else if(*eptr == 'Y') { wptr = wildcards[2]; }
        else if(*eptr == 'Z') { wptr = wildcards[3]; }

        while(*wptr != '\0') {
            if(*wptr == *sptr){
                wildcardmatch = TRUE;
                break;
            }
            wptr++;
        }
        if((wildcardmatch == TRUE) || (*eptr == *sptr)) {

            tnlargestpiece++;
            if(DEBUG_ALIGN) fprintf(stdout, "FOUND a match:
len=%d, epo(%d, %c), lib[%d].name=%s (%d, %c)\n",
                                tnlargestpiece, ecount, *eptr,
ilib, library[ilib].name, lcount, *sptr);

            if(tnlargestpiece > nlargestpiece) {
                nlargestpiece = tnlargestpiece;
                largestpiece_sptr = hold_this_place_sptr;
                largestpiece_lcount = hold_this_lcount;
                largestpiece_eptr = hold_this_place_eptr;
                largestpiece_ecount = hold_this_ecount;
            }
        }
    }
}

```

```

        if(DEBUG_ALIGN) fprintf(stdout, "FOUND a
largest piece: len=%d, epo(%d, %c), lib(%d, %c)\n",
                                nlargestpiece, largestpiece_ecount,
                                *largestpiece_eptr, largestpiece_lcount, *largestpiece_sptr);
    }
    sptr++;
    lcount++;
    eptr++;
    ecount++;
} else {
    tnlargestpiece = 0;
    sptr++;
    lcount++;
    hold_this_place_sptr = sptr;
    hold_this_lcount = lcount;
    /* NEW */
    eptr = hold_this_place_eptr;
    ecount = hold_this_ecount;
}
} else {
    tnlargestpiece = 0;
    break;
}
}
tnlargestpiece = 0;
eptr = hold_this_place_eptr + 1;
ecount = hold_this_ecount + 1;
hold_this_place_eptr = eptr;
hold_this_ecount = ecount;
if(DEBUG_ALIGN) {
    fprintf(stdout, "incrementing hold_this_place_eptr=%c,
hold_this_ecount=%d\n", *hold_this_place_eptr, hold_this_ecount);
}
}
}

```

```
if(DEBUG_ALIGN) fprintf(stdout,"ALIGN: largest piece match is %d monomers from
%s\n",nlargestpiece,library[ilib].name);
if(DEBUG_ALIGN) fprintf(stdout,"ALIGN: largestpiece_ecount=%d,
largestpiece_lcount=%d\n",
largestpiece_ecount,largestpiece_lcount);

if(nlargestpiece >= smallest_acceptable_piece) {
    if(DEBUG_ALIGN) fprintf(stdout,"ALIGN: incorporated\n");

    lcount = largestpiece_lcount;
    ecount = largestpiece_ecount;

    fprintf(stdout,"ALIGN_TARGET: ");
    for(ii=0; ii<largestpiece_ecount; ii++) {
        fprintf(stdout, " ");
    }

    while(ecount < (nlargestpiece + largestpiece_ecount)) {

        epothilone->alignedsequence[ecount] =
library[ilib].monomersequence[lcount];
        strcpy(epothilone->alignedPKSname[ecount],library[ilib].name);
        strcpy(epothilone->context[ecount],library[ilib].context[lcount]);
        epothilone->marked[ecount] = TRUE;
        if(ecount < (nlargestpiece + largestpiece_ecount - 1)) epothilone-
>boundarytoright[ecount] = FALSE;

        fprintf(stdout,"%c",library[ilib].monomersequence[lcount]);

        lcount++;
        ecount++;
    }

    for(ii=ecount; ii<elen; ii++) {
        fprintf(stdout, " ");
    }
}
```

```

    }
    fprintf(stdout,"%s\n",library[ilib].name);
}

return (nlargestpiece);

}/*maximal_adjacent_alignment_and_dump*/

/* -----
PURPOSE:
INPUT:
OUTPUT:
PROCEDURE:
*/

int  output_fresh_alignment(
LIB  *epotemp,
int  nboundary_cutoff)
{
int  account=0,ecount=0;
int  longest_segmentlen=0,current_segmentlen=0;
char *aptr,*eptr;
char boundary[MAXNAMELEN];

    eptr = epotemp->monomersequence;
    ecount = 0;
    epotemp->nboundary = 0;
    longest_segmentlen = 0;
    current_segmentlen = 0;
    strcpy(boundary,epotemp->alignedPKSname[ecount]);
    while(*eptr != '\0') {
        if(epotemp->boundarytoright[ecount] == TRUE) {

```

```

        epotemp->nboundary++;
        if(current_segmentlen > longest_segmentlen) {
            longest_segmentlen = current_segmentlen;
        }
        current_segmentlen = 0;
    }
    current_segmentlen++;
    ecount++;
    eptr++;
}
if(current_segmentlen > longest_segmentlen) {
    longest_segmentlen = current_segmentlen;
}

if(epotemp->nboundary > nboundary_cutoff) return 1;

eptr = epotemp->monomersequence;
ecount = 0;
fprintf(stdout, "HIT  ");
while(*eptr != '\0') {
    if(epotemp->alignedPKSname[ecount][0] == '\0') {
        if(epotemp->boundarytoright[ecount] == TRUE) {
            fprintf(stdout, "%c:TARG(%s)    ", *eptr, epotemp-
>context[ecount]);
        } else {
            fprintf(stdout, "%c:TARG(%s)    ", *eptr, epotemp-
>context[ecount]);
        }
    } else {
        if(epotemp->boundarytoright[ecount] == TRUE) {
            fprintf(stdout, "%c:%4s(%s) ", *eptr, epotemp-
>alignedPKSname[ecount], epotemp->context[ecount]);
        } else {
            fprintf(stdout, "%c:%4s(%s) ", *eptr, epotemp-
>alignedPKSname[ecount], epotemp->context[ecount]);
        }
    }
}

```

```

        }
    }
    ecount++;
    eptr++;
}

fprintf(stdout,"%d    %d",epotemp->nboundary,longest_segmentlen);
fprintf(stdout,"");

eptr = epotemp->monomersequence;
account = 0;
ecount = 0;
while(*eptr != '\0') {
    if(epotemp->boundarytoright[ecount] == TRUE) {
        fprintf(stdout,"%c|",epotemp->context[account][1]);
    } else {
        fprintf(stdout,"%c",epotemp->context[account][1]);
    }
    eptr++;
    account++;
    ecount++;
}

fprintf(stdout,"\n");
return 1;

}/*output_fresh_alignment*/

/* -----
*/
int  get_library(
char *libraryfile,
LIB  *library,
int  inter_modular_db_flag_off)
{

```



```
int    ii=0,jj=0,kk=0,lcount=0;
int    nlib=0;
char    *bufptr,buf[MAXBUF];
char    tmonomersequence[MAXNAMELEN];
char    *lptr,*tptr;
FILE    *libp;

    for(kk=0; kk < MAX_LIB_ENTRIES; kk++) {
        library[kk].recursion_tagged = FALSE;
        for(ii=0; ii<MAXNAMELEN; ii++) {
            library[kk].name[ii] = '\0';
            library[kk].monomersequence[ii] = '\0';
            library[kk].alignedsequence[ii] = '\0';

            for(jj=0; jj<MAXLEN; jj++) {
                library[kk].alignedPKSname[jj][ii] = '\0';
                library[kk].marked[jj] = FALSE;
                library[kk].context[jj][0] = '\0';
                library[kk].context[jj][1] = '\0';
                library[kk].context[jj][2] = '\0';
                library[kk].context[jj][3] = '\0';
            }
        }
        library[kk].nboundary = 0;
    }

/*    read in the library from PKS.lib */
if(NULL==(libp=fopen(libraryfile,"r"))) {
    fprintf(stdout,"TRY AGAIN; couldn't open library file: %s\n",libraryfile);
    nlib=0;
    exit();
}

nlib=0;
while(nlib < MAX_LIB_ENTRIES) {
```

```

        if(NULL==fgets(buf,sizeof(buf),libp)) break;
        bufptr = buf;
        if(*bufptr == '#') continue;
        if(*bufptr != '\n') {
/*
            sscanf(bufptr, "%s %s
%s", library[nlib].name, tmonomersequence, library[nlib].annotatedsequence);
*/

            lptr = library[nlib].name;
            while((*bufptr != ' ') && (*bufptr != '\0') && (*bufptr != '\n')){
                *lptr++ = *bufptr++;
            }
            *lptr = '\0';
            if((*bufptr != '\0') && (*bufptr != '\n')) bufptr++;

            lptr = library[nlib].monomersequence;
            while((*bufptr != ' ') && (*bufptr != '\0') && (*bufptr != '\n')){
/* -----
            This code specifically deletes inter-modular double bonds when the -d
option is set.
            ----- */
                if(inter_modular_db_flag_off == TRUE){
                    if(*bufptr != '=') {
                        *lptr++ = *bufptr++;
                    } else {
                        bufptr++;
                    }
                } else {
                    *lptr++ = *bufptr++;
                }
            }
            *lptr = '\0';
            if((*bufptr != '\0') && (*bufptr != '\n')) bufptr++;

```

```

        lptr = library[nlib].annotatedsequence;
        while((*bufptr != ' ') && (*bufptr != '\0') && (*bufptr != '\n')){
            *lptr++ = *bufptr++;
        }
        *lptr = '\0';
        if((*bufptr != '\0') && (*bufptr != '\n')) bufptr++;

        fprintf(stdout,"LIBRARY(%d) %s:
%s\n",nlib,library[nlib].name,library[nlib].monomersequence);
        fprintf(stdout,"LIBRARY(%d) %s:
%s\n",nlib,library[nlib].name,library[nlib].annotatedsequence);

        nlib++;
    }
}
fclose(libp);

for(kk=0; kk< nlib; kk++) {

    lcount = 0;
    lptr = library[kk].monomersequence;
    while(*lptr != '\0'){
        if(lcount == 0) {
            library[kk].context[lcount][0] = '-';
            library[kk].context[lcount][1] = *lptr;
            library[kk].context[lcount][2] = *(lptr + 1);
        } else {
            if(lcount == (strlen(library[kk].monomersequence) - 1)){
                library[kk].context[lcount][0] = *(lptr - 1);
                library[kk].context[lcount][1] = *lptr;
                library[kk].context[lcount][2] = '-';
            } else {
                library[kk].context[lcount][0] = *(lptr - 1);
                library[kk].context[lcount][1] = *lptr;

```

```

        library[kk].context[lcount][2] = *(lptr + 1);
    }
}
library[kk].context[lcount][3] = '\0';
lptr++;
lcount++;
}
fprintf(stdout, "LIBRARY(%d) %s:
%s\n", kk, library[kk].name, library[kk].monomersequence);
fprintf(stdout, "LIBRARY(%d) %s:
%s\n", kk, library[kk].name, library[kk].annotatedsequence);
for(jj=0; jj< strlen(library[kk].monomersequence); jj++) {
    fprintf(stdout, "(%s)", library[kk].context[jj]);
}
fprintf(stdout, "\n");
}

return nlib;
}/*get_library*/

```

```

/* -----
*/
int  dump_STARTER_align(
LIB  epothilone,
int  nwildcard,
char  wildcards[MAXWILD][MAXLEN]
)
{
int  elen=0;
int  ecount=0, hold_ecount=0;
int  wildcardmatch=FALSE;
char  *sptr, *eptr, *wptr;

    elen = strlen(epothilone.monomersequence);

```

```
eptr = epothilone.monomersequence;
fprintf(stdout, "ALIGN_TARGET: ");
while(*eptr != '\0') {
    fprintf(stdout, "%c", *eptr);
    eptr++;
}
fprintf(stdout, "\n");

ecount=0;
eptr = epothilone.monomersequence;
sptr = epothilone.alignedsequence;
fprintf(stdout, "ALIGN_TARGET: ");
while(*eptr != '\0') {
    wildcardmatch = FALSE;
    wptr = "";
    if(*eptr == 'X') { wptr = wildcards[0]; }
    else if(*eptr == 'Y') { wptr = wildcards[1]; }
    else if(*eptr == 'Z') { wptr = wildcards[2]; }

    while(*wptr != '\0') {
        if(*wptr == *sptr){
            wildcardmatch = TRUE;
            break;
        }
        wptr++;
    }
    if(wildcardmatch == TRUE) {
        fprintf(stdout, ".");
    } else {
        if(*eptr == *sptr) {
            fprintf(stdout, "|");
        } else {
            fprintf(stdout, " ");
        }
    }
}
```

```

        eptr++;
        sptr++;
    }
    fprintf(stdout, "\n");
    fprintf(stdout, "ALIGN_TARGET: ");
    eptr = epothilone.monomersequence;
    ecount=0;
    while(*eptr != '\0') {
        if(epothilone.alignedsequence[ecount] == '\0') {
            fprintf(stdout, " ");
        } else {
            fprintf(stdout, "%c", epothilone.alignedsequence[ecount]);
            hold_ecount = ecount;
        }
        eptr++;
        ecount++;
    }
    fprintf(stdout, "      %s\n", epothilone.alignedPKSname[hold_ecount]);
    fprintf(stdout, "STARTER_ALIGN:\n");

} /*dump_STARTER_align*/

/* -----
*/
int  reset_epotemp(
LIB  *epotemp,
LIB  epothilone)
{
int  jj=0, elen=0;

    elen = strlen(epothilone.monomersequence);

    strcpy(epotemp->name, epothilone.name);
    strcpy(epotemp->monomersequence, epothilone.monomersequence);
    strcpy(epotemp->alignedsequence, epothilone.alignedsequence);

```

```
epotemp->nboundary = 0;
for(jj=0; jj< elen; jj++) {
    strcpy(epotemp->alignedPKSname[jj],epothilone.alignedPKSname[jj]);
    epotemp->marked[jj] = epothilone.marked[jj];
    strcpy(epotemp->context[jj],epothilone.context[jj]);
    epotemp->boundarytoright[jj] = epothilone.boundarytoright[jj];
}

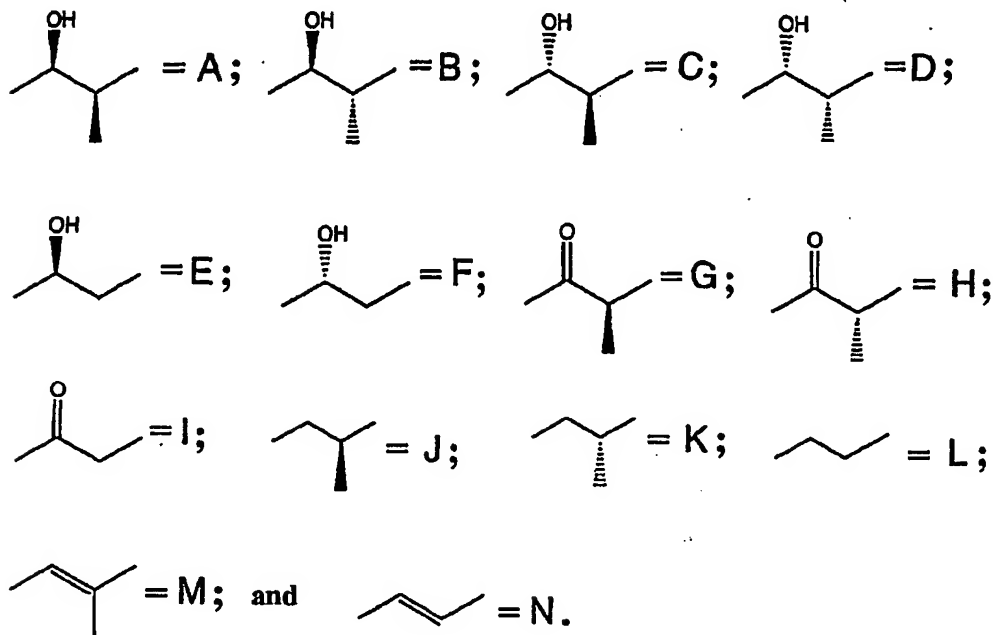
}/*reset_epotemp*/
```

Thus, the present invention provides a useful means to generate new PKS genes and corresponding enzymes to produce polyketides. The invention having now been described by way of written description and examples, those of skill in the art will recognize that the invention can be practiced in a variety of embodiments and that the foregoing description and examples are for purposes of illustration and not limitation of the following claims.

WHAT IS CLAIMED IS:

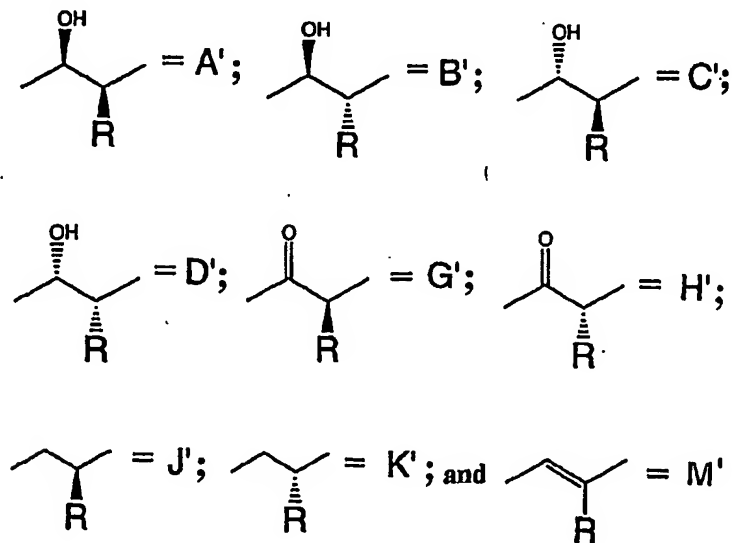
1. A method for representing the structure of a polyketide produced by a modular polyketide synthase, said method comprising the steps of:
 - (a) defining a set of monomer units of which said polyketide is composed,
 - (b) assigning an alphanumeric symbol or symbols to each different monomer unit in said set,
 - (c) identifying one or more monomers in said set that is present in said polyketide, and
 - (d) composing a string of said symbols ordered in a manner reflecting the order in which said monomers occurs in said polyketide, wherein said string of symbols represents the structure of said polyketide.
2. The method of claim 1, wherein said monomer set comprises two-carbon unit monomers, wherein a first carbon of said unit is substituted with hydrogen or methyl, and a second carbon of said unit is substituted with oxygen, hydroxy, or hydrogen, and said two carbon unit comprises either a single or a double bond between said first and second carbons.
3. The method of claim 2, wherein said monomer set additionally comprises one or more members selected from the group consisting of two carbon unit monomers in which said first carbon is substituted with hydroxy, methoxy, or ethyl; a moiety corresponding to an amino acid or amino acid derivative incorporated into a PKS by a non-ribosomal peptide synthase; a moiety corresponding to a structure incorporated into a polyketide by an AMP ligase or a CoA ligase; and a moiety corresponding to a structure corresponding to a structure in a polyketide after modification by a polyketide modification enzyme.

4. The method of claim 2 wherein the set of monomer unit and corresponding symbol comprises:

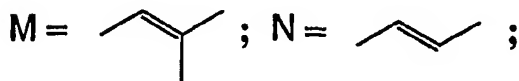
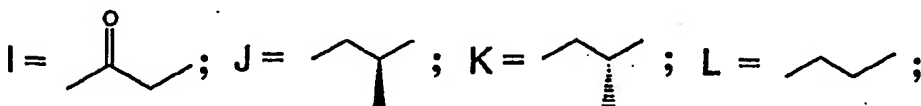
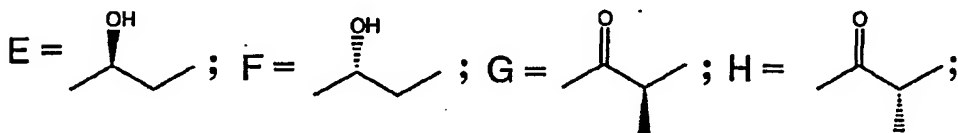
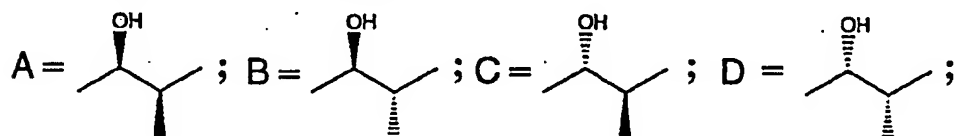


5. The method of claim 4 wherein the set of monomer unit further comprises a miscellaneous monomer that is assigned the symbol Q.

6. The method of claim 4 wherein the set of monomer unit and corresponding symbol further comprises

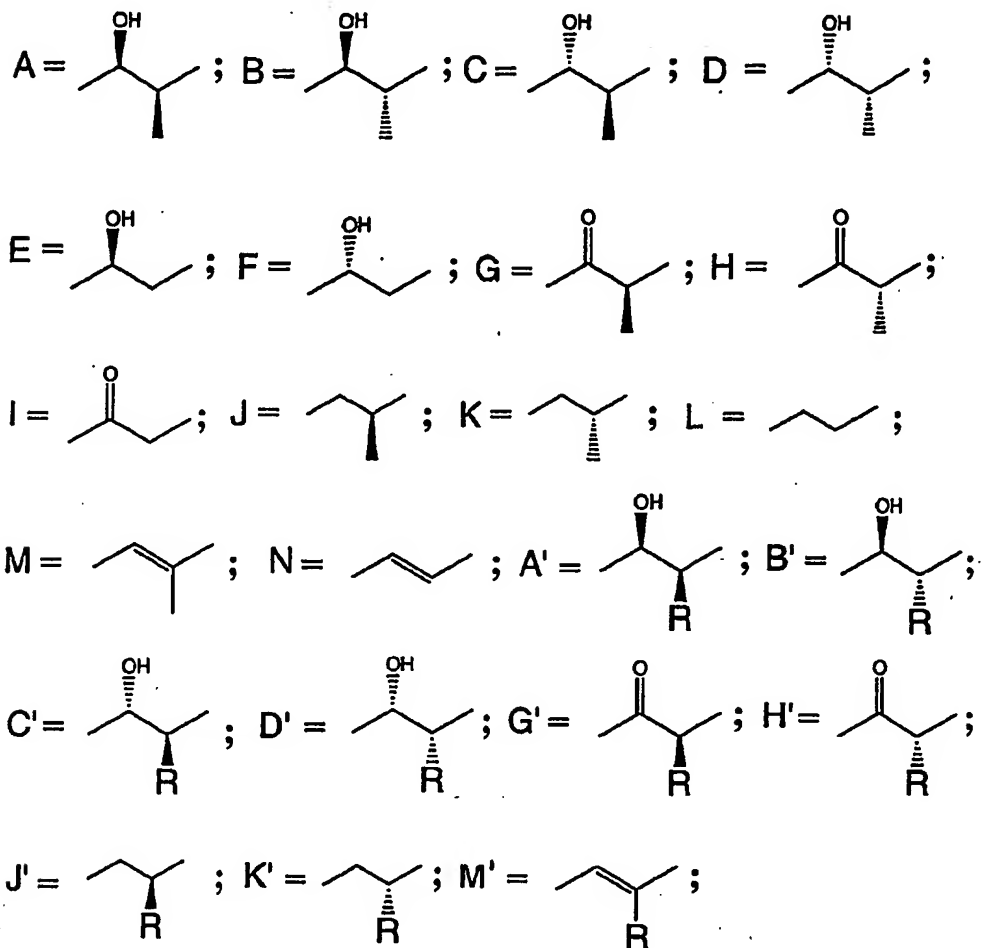


7. A database of polyketides, in which each said member is represented by a string of alpha-numeric symbols, wherein said symbols represent structural subunits of said polyketide, and said string represents the order in which such subunits occur in said polyketide.
8. The database of claim 7 that includes at least 100 different polyketides.
9. The database of claim 7 wherein each said member is represented by a CHUCKLES string.
10. The database of claim 7 wherein each said member is represented by an annotated CHUCKLES string.
11. The database of claim 7 wherein the symbol and its corresponding structural subunit are selected from the group consisting of



and Q for a miscellaneous monomer.

12. The database of claim 7 wherein the symbol and its corresponding structural subunit are selected from the group consisting of



and Q for a miscellaneous monomer.

13. A database of polyketides, in which each said member is represented by a linearized representation of said polyketide.

14. A method of designing a PKS gene capable of producing a desired polyketide, which method comprises:

- (a) defining a string of alphanumeric symbols representing the structure of said polyketide,
- (b) comparing said string to a database of strings of alphanumeric symbols representing polyketides produced by PKS genes,

(c) identifying common elements in said string representing the structure of said polyketide with elements in said strings in said database, and

(d) generating one or more new strings from elements identified in step (b) that match said string representing the structure of said polyketide, wherein said new string defines a PKS gene capable of producing said polyketide.

15. The method of claim 14, wherein all possible PKS genes encoding a desired polyketide from said database are generated and displayed.

16. The method of claim 14, wherein said new strings generated in step (d) are rated and displayed in an order based on one or more parameters.

17. The method of claim 16, wherein said parameters are selected from the group consisting of number of non-native module interfaces and number of non-native protein interfaces.

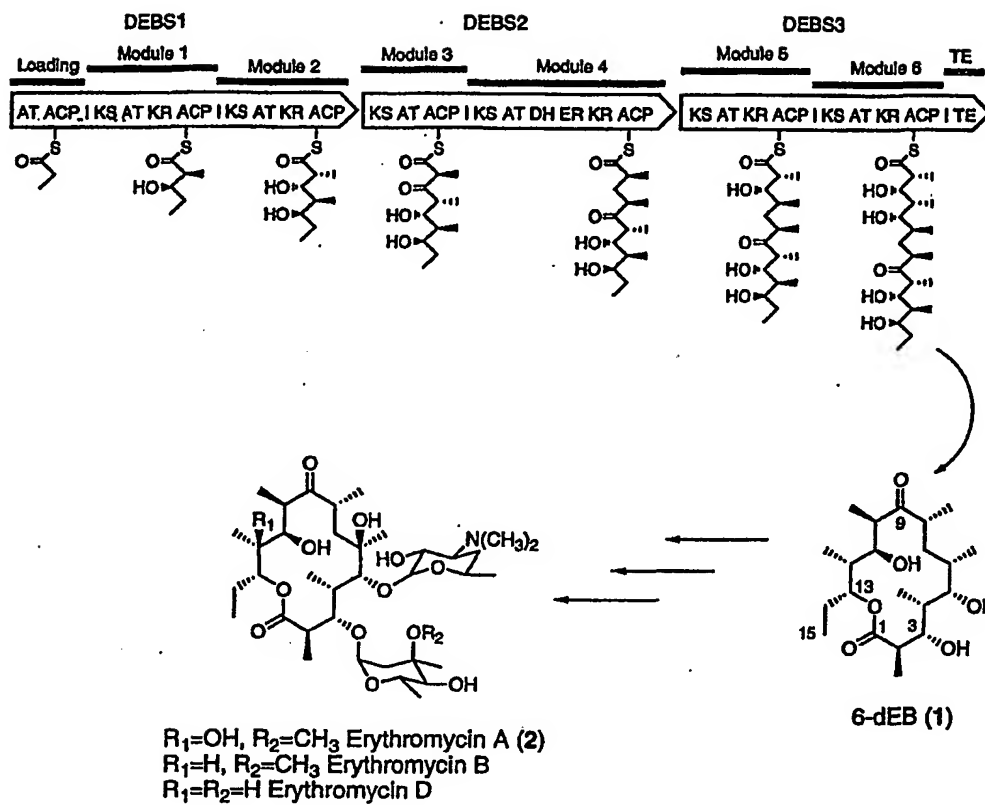
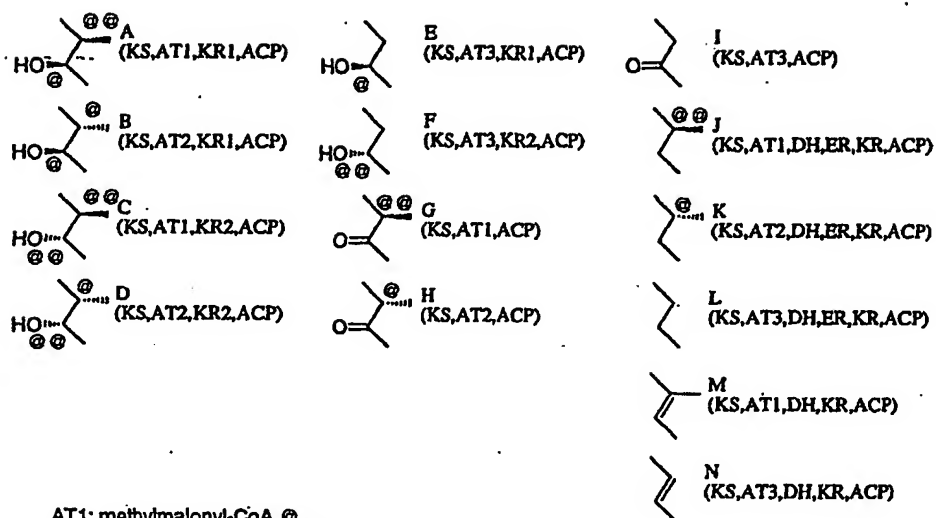


Figure 1

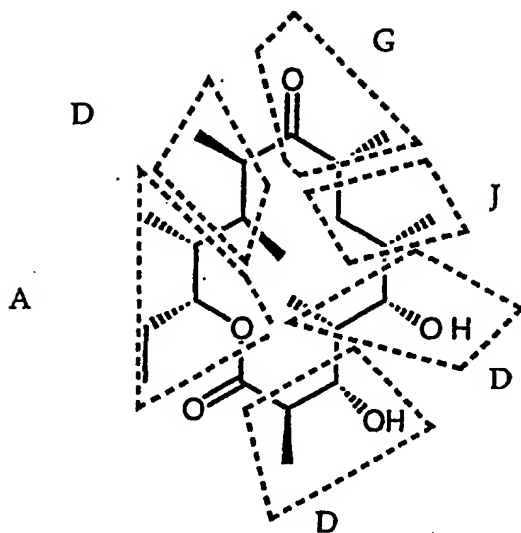


AT1: methylmalonyl-CoA @
 AT2: methylmalonyl-CoA @@
 AT3: malonyl-CoA

KR1: hydroxyl @
 KR2: hydroxyl @@

@@ = counter clockwise
 @ = clockwise

Figure 2



CHUCKLES: ADGJDD

SMILES:

```
C1(=O)-[C@H](C)[C@@H](OH)-[C@@H](C)[C@@H](OH)-  
[C@@H](C)C-[C@@H](C)C(=O)-[C@H](C)[C@@H](C)-  
[C@@H](C)[C@@H](CC)O1
```

Figure 3

Figure 4A

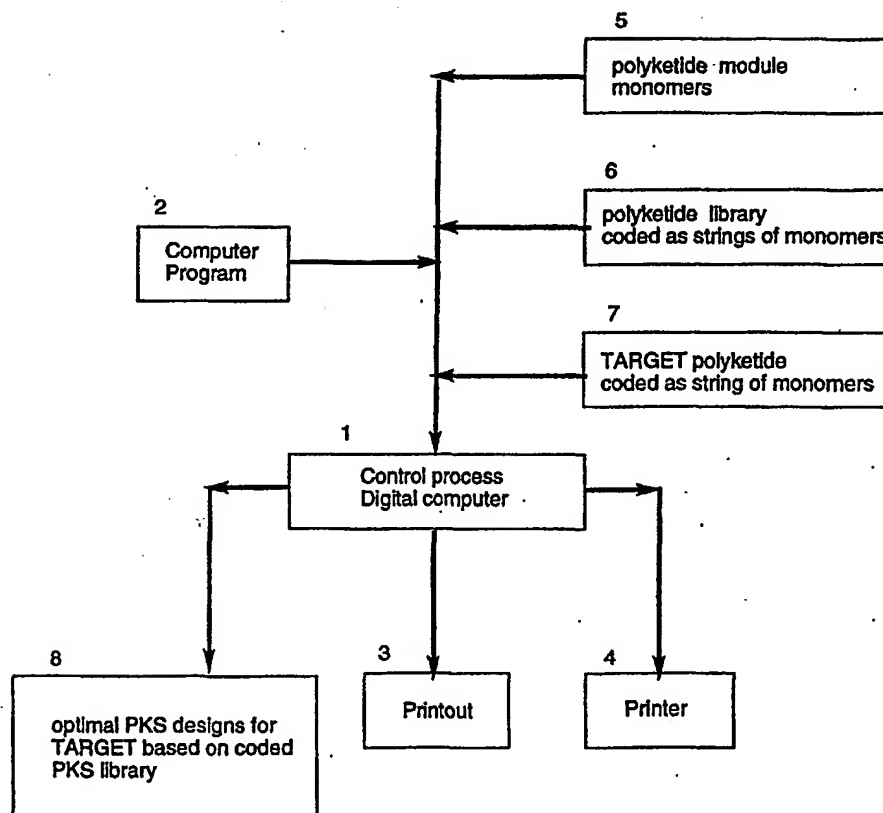


Figure 4B

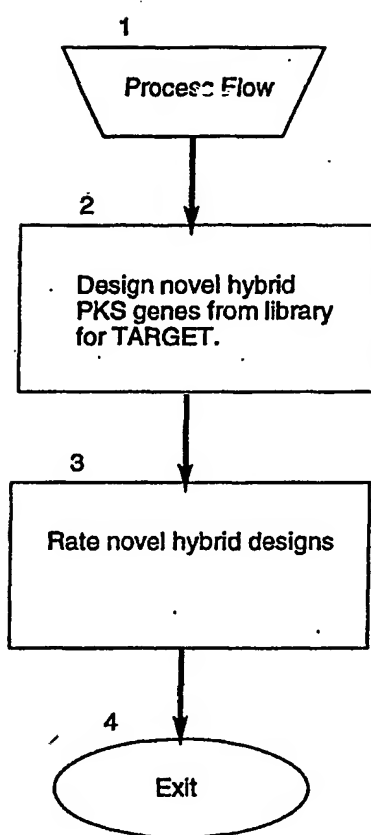
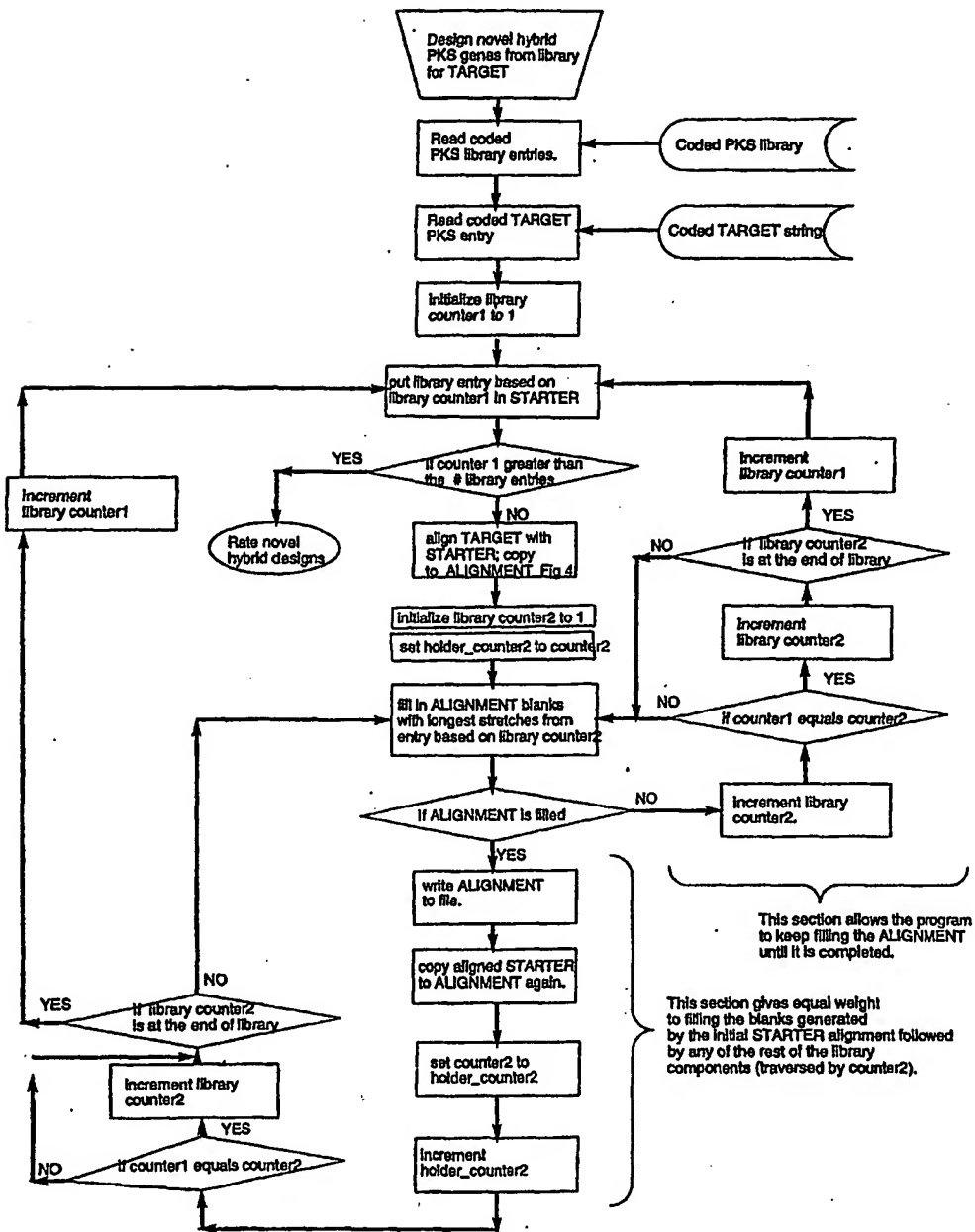


Figure 4C



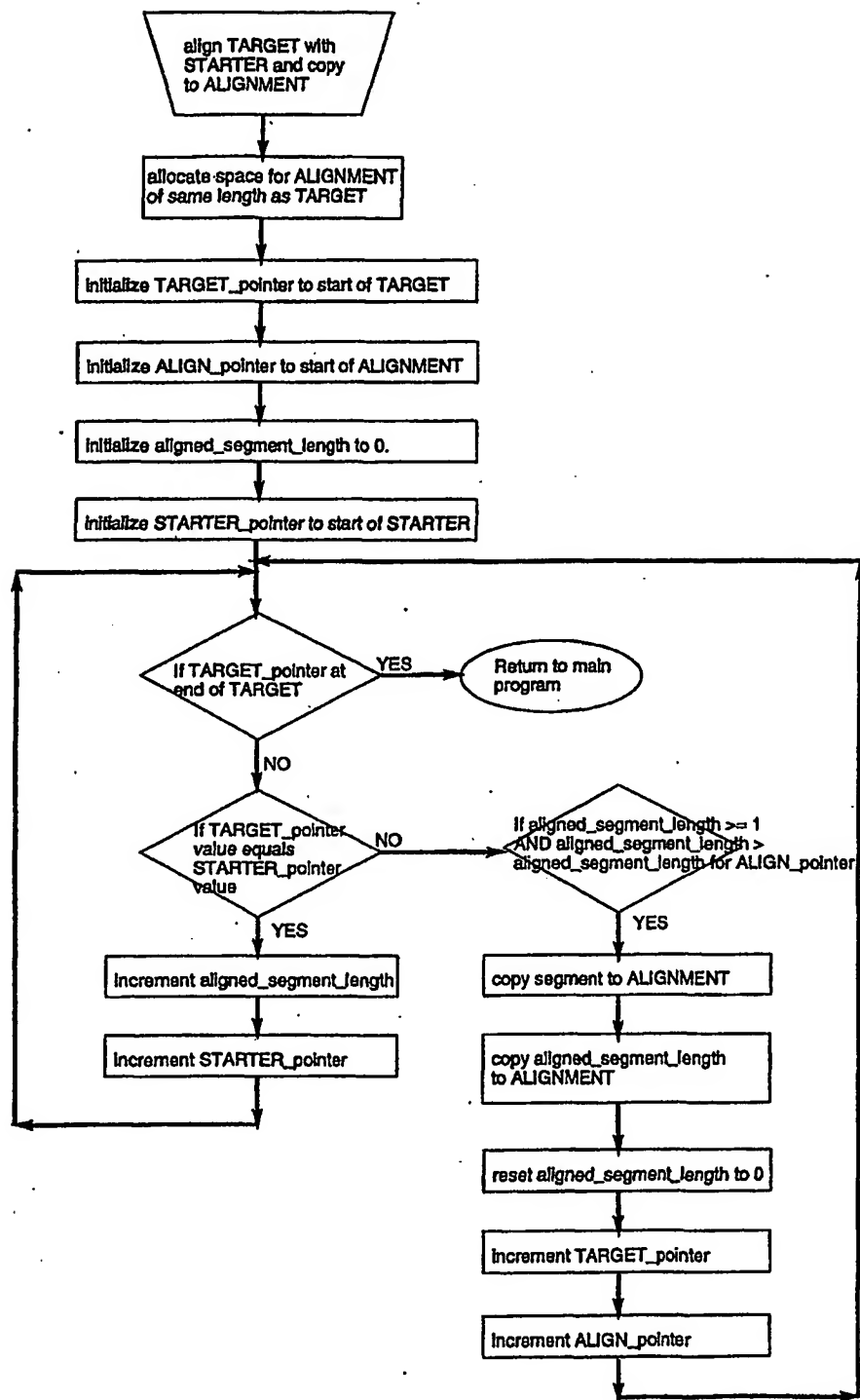


Figure 4D

Figure 4E

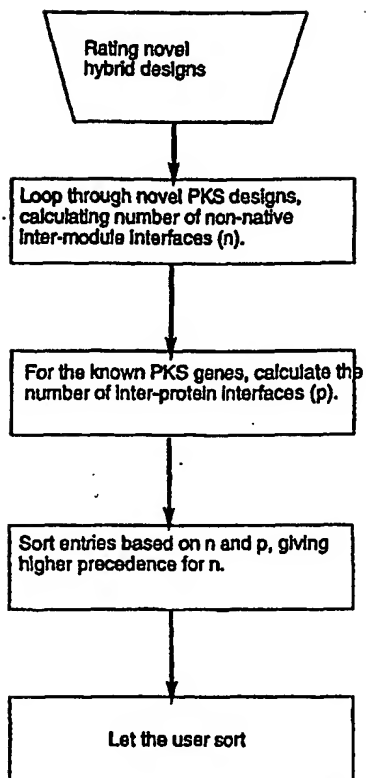
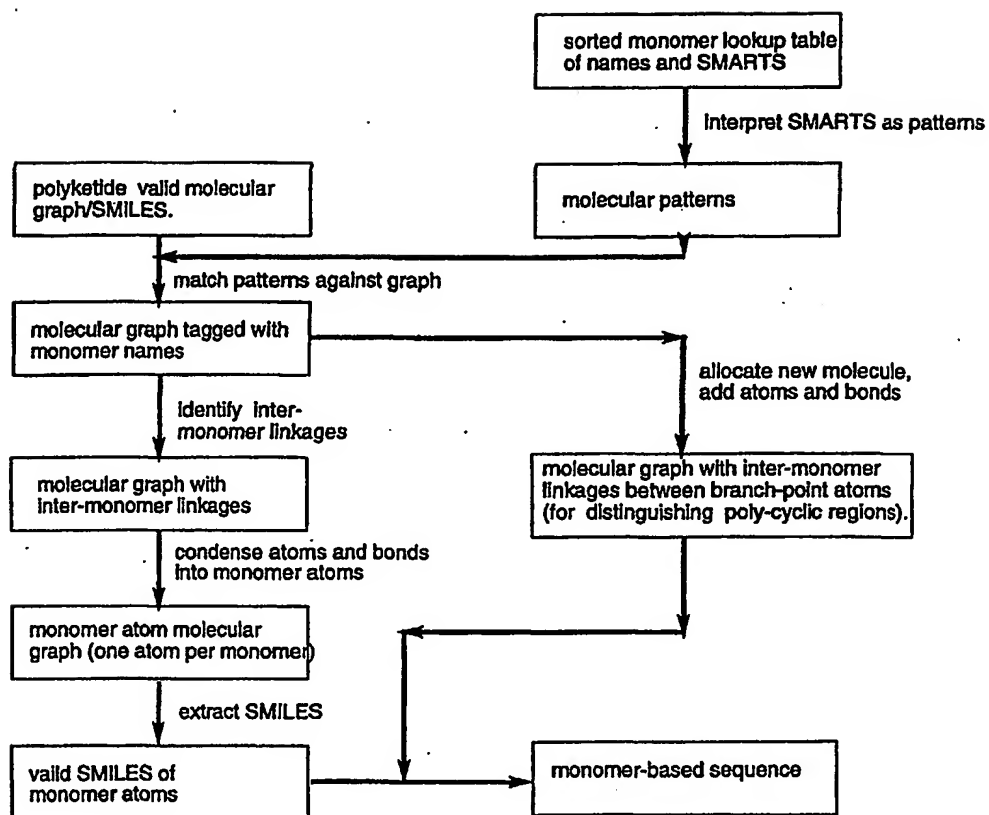


Figure 5



(12) INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
6 December 2001 (06.12.2001)

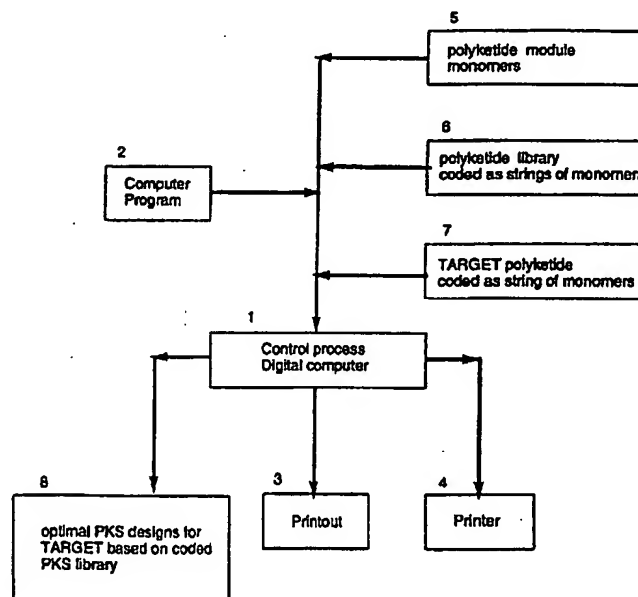
PCT

(10) International Publication Number
WO 01/092991 A3

- (51) International Patent Classification⁷: **C12Q 1/48, G06F 19/00**
- (21) International Application Number: **PCT/US01/17352**
- (22) International Filing Date: **29 May 2001 (29.05.2001)**
- (25) Filing Language: **English**
- (26) Publication Language: **English**
- (30) Priority Data:
60/207,331 30 May 2000 (30.05.2000) US
60/237,382 4 October 2000 (04.10.2000) US
- (71) Applicant: **KOSAN BIOSCIENCES, INC. [US/US];**
3832 Bay Center Drive, Hayward, CA 94545 (US).
- (72) Inventors: **KHOSLA, Chaitan; 740 La Para Avenue, Palo Alto, CA 94306 (US). REID, Ralph, C.; 600 Galerita Way, San Rafael, CA 94903 (US). SANTI, Daniel, V.; 211 Belgrave Avenue, San Francisco, CA 94417 (US). SIANI, Michael, A.; 341 Day Street, San Francisco, CA 94131 (US).**
- (74) Agents: **FAVORITO, Carolyn, A. et al.; Morrison & Forster LLP, Suite 500, 3811 Valley Centre Drive, San Diego, CA 92130-2332 (US).**
- (81) Designated States (*national*): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, TZ, UA, UG, UZ, VN, YU, ZA, ZW.
- (84) Designated States (*regional*): ARIPO patent (GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE, TR), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).
- Published:
— with international search report

[Continued on next page]

(54) Title: DESIGN OF POLYKETIDE SYNTHASE GENES



(57) Abstract: Methods for the computational analysis of polyketides and the computer-assisted design of PKS genes are facilitated by representing the structure of a polyketide and/or a PKS gene that encodes the PKS that produces the polyketide by alphanumeric symbols that facilitates computer assisted analysis. A database of polyketides and corresponding PKS genes that can be rapidly searched and information extracted for a variety of applications, including the design and specification of PKS genes via the recombining of modules or portions of modules or sets of modules from already known and variable PKS genes.

WO 01/092991 A3



(88) Date of publication of the international search report:
8 August 2002

For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

INTERNATIONAL SEARCH REPORT

onal Application No
PCT/US 01/17352

A. CLASSIFICATION OF SUBJECT MATTER

IPC 7 C12Q1/48 G06F19/00

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

IPC 7 C12Q G06F

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

EPO-Internal, WPI Data, PAJ, BIOSIS, MEDLINE, EMBASE, SCISEARCH, CHEM ABS Data

C. DOCUMENTS CONSIDERED TO BE RELEVANT

| Category * | Citation of document, with indication, where appropriate, of the relevant passages | Relevant to claim No. |
|------------|---|-----------------------|
| Y | SIANI, M.A. ET AL.: "Chuckles: A method for representing and searching peptide and peptoid sequences on both monomer and atomic levels" J. CHEM. INF. COMPUT. SCI., vol. 54, 1994, pages 588-593, XP001062074 cited in the application the whole document | 1-17 |
| Y | XUE Q ET AL: "MULTIPLASMID APPROACH TO PREPARING LARGE LIBRARIES OF POLYKETIDES" PROCEEDINGS OF THE NATIONAL ACADEMY OF SCIENCES OF USA, NATIONAL ACADEMY OF SCIENCE. WASHINGTON, US, vol. 96, no. 21, 1999, pages 11740-11745, XP000910245 ISSN: 0027-8424 the whole document | 1-17 |



Further documents are listed in the continuation of box C.



Patent family members are listed in annex.

* Special categories of cited documents:

- *A* document defining the general state of the art which is not considered to be of particular relevance
- *E* earlier document but published on or after the international filing date
- *L* document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)
- *O* document referring to an oral disclosure, use, exhibition or other means
- *P* document published prior to the international filing date but later than the priority date claimed

T later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

X document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

Y document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.

G document member of the same patent family

Date of the actual completion of the international search

22 April 2002

Date of mailing of the international search report

06/05/2002

Name and mailing address of the ISA

European Patent Office, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel. (+31-70) 340-2040, Tx. 31 651 epo nl,
Fax: (+31-70) 340-3016

Authorized officer

Moreno de Vega, C.

INTERNATIONAL SEARCH REPORT

International Application No
PCT/US 01/17352

C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT

| Category * | Citation of document, with indication, where appropriate, of the relevant passages | Relevant to claim No. |
|------------|--|-----------------------|
| Y | KHOSLA C ET AL: "Generation of polyketide libraries via combinatorial biosynthesis" TRENDS IN BIOTECHNOLOGY, ELSEVIER PUBLICATIONS, CAMBRIDGE, GB, vol. 14, no. 9, 1 September 1996 (1996-09-01), pages 335-341, XP004035738 ISSN: 0167-7799 the whole document | 1-17 |
| P,X | SIANI, M.A. ET AL: "Development and screening of a polyketide virtual library for drug leads against a motilide pharmacophore" JOURNAL OF MOLECULAR GRAPHICS AND MODELLING, vol. 18, no. 4-5, August 2000 (2000-08), pages 497-511, XP001068499 the whole document | 1-17 |
| Y | WO 98 27203 A (KOSAN BIOSCIENCES) 25 June 1998 (1998-06-25) the whole document | 1-17 |
| Y | WO 00 24907 A (KOSAN BIOSCIENCES) 4 May 2000 (2000-05-04) the whole document | 1-17 |

INTERNATIONAL SEARCH REPORT

Initial Application No
PCT/US 01/17352

| Patent document cited in search report | | Publication date | Patent family member(s) | Publication date |
|---|---|---------------------|----------------------------|---------------------|
| WO 9827203 | A | 25-06-1998 | AU 734325 B2 | 07-06-2001 |
| | | | AU 5701098 A | 15-07-1998 |
| | | | EP 0948613 A1 | 13-10-1999 |
| | | | JP 2001510993 T | 07-08-2001 |
| | | | WO 9827203 A1 | 25-06-1998 |
| | | | US 6258566 B1 | 10-07-2001 |
| | | | US 6033883 A | 07-03-2000 |
| <hr/> | | | | |
| WO 0024907 | A | 04-05-2000 | AU 1447700 A | 15-05-2000 |
| | | | EP 1124968 A2 | 22-08-2001 |
| | | | WO 0024907 A2 | 04-05-2000 |
| <hr/> | | | | |